

ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
ΣΧΟΛΗ ΘΕΤΙΚΩΝ ΕΠΙΣΤΗΜΩΝ
ΤΜΗΜΑ ΕΠΙΣΤΗΜΗΣ ΥΠΟΛΟΓΙΣΤΩΝ

**Σχεδιασμός και υλοποίηση ενός Ανα-Ροή
Διαχειριστή ουρών για ένα μεταγωγέα τύπου ATM με
χρήση τεχνολογίας FPGA**

Δημήτριος Σ. Καψάλης

Μεταπτυχιακή Εργασία

Ηράκλειο, Φεβρουάριος 2002

ΠΑΝΕΠΙΣΤΗΜΙΟ ΚΡΗΤΗΣ
ΣΧΟΛΗ ΘΕΤΙΚΩΝ ΕΠΙΣΤΗΜΩΝ
ΤΜΗΜΑ ΕΠΙΣΤΗΜΗΣ ΥΠΟΛΟΓΙΣΤΩΝ

Σχεδιασμός και υλοποίηση ενός Ανα-Ροή Διαχειριστή ουρών για ένα μεταγωγέα τύπου ATM με χρήση τεχνολογίας FPGA

Εργασία που υποβλήθηκε από τον
ΔΗΜΗΤΡΙΟ Σ. ΚΑΨΑΛΗ
ως μερική εκπλήρωση των απαιτήσεων
για την απόκτηση
ΜΕΤΑΠΤΥΧΙΑΚΟΥ ΔΙΠΛΩΜΑΤΟΣ ΕΙΔΙΚΕΥΣΗΣ

Συγγραφέας:

Δημήτριος Σ. Καψάλης
Τμήμα Επιστήμης Υπολογιστών
Πανεπιστήμιο Κρήτης

Εξεταστική Επιτροπή:

Μανώλης Κατεβαίνης, Καθηγητής
Επόπτης

Απόστολος Τραγανίτης, Αναπληρωτής Καθηγητής
Μέλος

Ευάγγελος Μαρκάτος, Αναπληρωτής Καθηγητής
Μέλος

Δεκτή:

Πάνος Κωνσταντόπουλος, Καθηγητής
Πρόεδρος Επιτροπής Μεταπτυχιακών Σπουδών

Φεβρουάριος 2002

Σχεδιασμός και υλοποίηση ενός Ανα-Ροή Διαχειριστή ουρών για ένα μεταγωγέα τύπου ATM με χρήση τεχνολογίας FPGA

Δημήτριος Σ. Καψάλης

Μεταπτυχιακή Εργασία

Τμήμα Επιστήμης Υπολογιστών
Πανεπιστήμιο Κρήτης

Περίληψη

Οι προχωρημένοι Μεταγωγείς και Δρομολογητές στηρίζονται κυρίως στην τεχνολογία Δυναμικής RAM για την παροχή μεγάλου, χαμηλού κόστους χώρου, που είναι απαραίτητος λόγω της εκρηκτικότητας της Διαδικτυακής κίνησης. Ποιότητα Υπηρεσίας (Quality of Service) είναι επίσης επιθυμητή. Κατά συνέπεια η Ανά-Ροή Αποθήκευση σε Ουρές (Per-Flow Queueing) συχνά υλοποιείται. Μελετάμε τον σχεδιασμό ενός Διαχειριστή Ουρών που υποστηρίζει Ανά-Ροή Αποθήκευση σε ουρές χιλιάδες ροών κίνησης τύπου ABR για ένα μεταγωγέα ATM. Ένα μεγάλο ολοκληρωμένο τύπου FPGA χρησιμοποιείται για γρήγορη ανάπτυξη και εκτενείς δοκιμές πάνω στην πλατφόρμα. Προς αποφυγή της χρήσης ολοκληρωμένων μνήμης τύπου SRAM και μείωση της χρήσης pin και συρμάτων, μόνο μία μονάδα μνήμης τύπου SDRAM DIMM χρησιμοποιείται για την αποθήκευση κελιών και τη διατήρηση δεικτών. Προτιμήσαμε τη Δυναμική Παραχώρηση Μνήμης (Dynamic Memory Allocation), προκειμένου να αντιμετωπιστούν οι ροές υψηλής κίνησης. Προπαραχώρηση αποθηκευτικών χώρων (Buffer Preallocation) και Παράκαμψη Λίστας Ελευθέρων χώρων (Free List bypassing) χρησιμοποιήθηκαν για την μείωση των προσπελάσεων μνήμης και την αύξηση της αποθηκευτικής διαμεταγωγής. Αυτές οι τεχνικές αποδεικνύονται απαραίτητες για την ικανοποίηση των αναγκών αποθήκευσης του Μεταγωγέα. Χαρακτηριστικά Ελέγχου Ροής τύπου ATM (ATM Flow Control), όπως Μαρκάρισμα τυπου EFCI (EFCI Marking) και Μαρκάρισμα RM σχετικού ρυθμού (RM Relative rate marking) παρέχεται για κάθε υποστηριζόμενη ροή. Χρησιμοποιήσαμε το συνθέσιμο υποσύνολο της γλώσσας περιγραφής υλικού Verilog για προσομοίωση και σχεδιασμό της αρχιτεκτονικής, αντί της ALTERA AHDL, για συμβατότητα μεταξύ διαφορετικών πλατφορμών. Το εργαλείο ALTERA MaxPlusII χρησιμοποιήθηκε για σύνθεση και προγραμματισμό της FPGA. Πετύχαμε 35 MHz συχνότητας ρολογιού που μεταφράζεται σε 800 Mbps μέγιστη συνδιασμένη, εισερχόμενη και εξερχόμενη διαμεταγωγή για τον Διαχειριστή Ουρών καθώς και μια πολυπλοκότητα 2500 λογικών στοιχείων FPGA (FPGA Logic Elements) και 2000 SRAM bit για 64 χιλιάδες ροές.

Επόπτης

Μανώλης Κατεβαίνης
Καθηγητής
Τμήμα Επιστήμης Υπολογιστών
Πανεπιστήμιο Κρήτης

Design and Implementation of a Per-Flow Queue Manager for an ATM Switch using FPGA technology

Dimitrios S. Kapsalis

Master of Science Thesis

Department of Computer Science
University of Crete
Greece

Abstract

Advanced Switches and routers rely mostly on Dynamic RAM technology for providing large, low-cost buffer space needed due to the burstiness of Internet traffic. Quality of Service is also desirable, therefore, per flow queueing of traffic is often implemented.

We designed and implemented a queue manager that supports per flow queueing of thousands of flows of ABR traffic for an ATM Switch. A large FPGA chip was used for fast development and extensive on-board testing. To avoid SRAM chip usage and lower the pin and trace count, a single SDRAM DIMM is used for storing both cells and pointers. We implemented dynamic memory allocation.

Buffer preallocation and free list bypassing were used to reduce memory accesses and increase buffer bandwidth. These techniques proved essential for satisfying the switch buffer requirements. ATM Flow Control features such as EFCI and RM Relative Rate Marking has been provided for each supported flow. We used synthesizable Verilog for simulation and of the architecture instead of ALTERA AHDL, so as to achieve cross-platform compatibility.

The ALTERA MaxPlus II tool has been used for synthesis and FPGA programming. We achieved a clock frequency of 35 MHz; this translates to a peak of 800 Mbps of combined incoming and outgoing throughput for the Queue Manager; the queue manager occupies 2500 FPGA Logic Elements and 2000 SRAM bits for 64K flows.

Advisor

Manolis Katevenis
Professor
Department of Computer Science

University of Crete
Greece

Ευχαριστίες

Η εργασία αυτή αποτελεί το επιστέγασμα των σπουδών μου στο Πανεπιστήμιο Κρήτης. Για την επίτευξη αυτού του στόχου είμαι βαθιά υπόχρεος σε όλους τους ανθρώπους που με επιρέασαν και μου ενέπνευσαν την επιθυμία για σπουδή και δημιουργία. Σ' αυτό το σημείο θα ήθελά να ευχαριστήσω τους ανθρώπους που είχαν κύρια συμβολή στην ολοκλήρωση αυτής της εργασίας.

Τον επόπτη και ακαδημαϊκό σύμβουλο, καθηγητή Μανώλη Κατεβαΐνη, για την καθοδήγηση και την υποστήριξη καθόλη την διάρκεια των μεταπτυχιακών σπουδών μου, αλλά και την συνεισφορά του στην αναζωογόνηση της πίστης μου στην επιστημονική γνώση και έρευνα ως τρόπο και στάση ζωής. Τον αρχικό επόπτη και ακαδημαϊκό αναπληρωτή καθηγητή σύμβουλο Δημήτρη Σερπάνο για την καθοδήγηση αλλά και την εμπιστοσύνη και την συνεισφορά του στην έναρξη των μεταπτυχιακών μου σπουδών. Τους καθηγητές του τμήματος κύριους Απόστολο Τραγανίτη και Βαγγέλη Μαρκάτο για την συμμετοχή τους στην εξεταστική επιτροπή της εργασίας.

Το Εργαστήριο Αρχιτεκτονικής και VLSI του Ιδρύματος Τεχνολογίας και Έρευνας (CARV-ICS) για την ευχάριστη συνεργασία και την πολυτιμη εμπειρία. Κυρίως τον Διονύση Πνευματικάτο, τον Γιώργο Καλοκαιρινό, και τον Μιχάλη Λιγεράκη για την καθοδήγηση και την πολύτιμη κριτική, για το ποιοτικό περιβάλλον εργασίας και τις όμορφες στιγμές που θα μου μείνουν αξέχαστες. Τον Χρήστο Λόλα και τον Γιώργο Παπαδάκη για την συνεργασία το κέφι και πάνω απ' όλα την αλληλεγγύη σε όλη το εύρος της παράλληλης πορείας μας που μερικές φορές φαινόταν ατελείωτη.

Το προσωπικό του Τμήματος Επιστήμης Υπολογιστών: της κυρίες Ρένα Καλαϊτζάκη και Μαρία Σταυρακάκη για την άμεση βοήθεια και την υποστήριξη στην εκπλήρωση των υποχρεώσεών μου.

Την Γενική Γραμματεία Έρευνας και Τεχνολογίας που χρηματοδότησε το έργο «ΔΙΠΟΛΟ». Το Ίδρυμα Τεχνολογίας και Έρευνας και το Πανεπιστήμιο Κρήτης που υποστήριξαν οικονομικά τόσο εμένα και ένα πλήθος συναδέλφων μου που εργστήκαμε για την υλοποίηση του έργου. Το προσωπικό του Ινστιτούτου Πληροφορικής και κυρίως την Μαρία Πρεβελιανάκη για την υποστήριξη και την πολύτιμη φιλία.

Τους συνάδελφούς μου στην ISD, τον Κώστα Παπαδά, τον Σπύρο Λυμπέρη, για την υποστήριξη, τις γνώσεις και την ευκαιρία για ένα νέο ξεκίνημα.

Τους φίλους μου, σε όλα τα χρόνια των σπουδών μου που τα γεμίσαν με το παραπάνω. Τον Μανώλη, την Κατερίνα, τον Ξενοφόντα και πολλούς άλλους....

Την Αργυρώ που ήταν τόσο μακριά, αλλά τόσο κοντά.

Πάνω απ' όλα στους γονείς μου, Στέφανο και Βάσιλική και την αδελφή μου Σοφία. Σ' αυτούς οφείλω την ύπαρξή μου και την υπόστασή μου. Ελπίζω να φανώ αντάξιος των θυσιών τους και της αγάπης τους.

Πίνακας Περιεχομένων

ΕΥΧΑΡΙΣΤΙΕΣ	9
ΠΙΝΑΚΑΣ ΠΕΡΙΕΧΟΜΕΝΩΝ	11
ΛΙΣΤΑ ΣΧΗΜΑΤΩΝ	14
ΛΙΣΤΑ ΠΙΝΑΚΩΝ	15
1 ΕΙΣΑΓΩΓΗ	17
1.1 Κίνητρα	17
1.2 Περιεχόμενα αυτής της εργασίας στο πλαίσιο του προγράμματος ΔΙΠΟΛΟ.	19
1.3 Η αρχιτεκτονική του μεταγωγέα ΔΙΠΟΛΟ	20
1.4 Αρχιτεκτονική του εξυπηρετητή ABR κίνησης	21
1.4.1 Το στοιχείο μεταγωγής	22
1.4.2 Εσωτερικός επεξεργαστής ελέγχου	23
1.4.3 Η μονάδα εξυπηρετή κίνησης ABR	23
1.4.4 Η μονάδα μνήμης SDRAM (256 MB) [19]	24
1.4.5 Μηχανισμός Flow Control Κάρτας Εξυπηρετητή ABR [23]	24
1.5 Η αρχιτεκτονική της Μονάδας Εξυπηρετητή ABR (ABRSU)	26
1.5.1 Η Διεπαφή Επεξεργαστή (CPU Interface)	27
1.5.2 Μονάδα διεπαφής CPU – διεπαφή με Διαχειριστή Ουρών	28
1.5.3 Μονάδα διεπαφής CPU - διεπαφή με Cell Scheduler	29
1.5.4 Η υπομονάδα Cell Demultiplexor	30
1.5.5 Διεπαφές UTOPIA	30
1.5.6 Ο Προγραμματιστής Cells (Cell Scheduler) [23]	31
1.5.7 Ο Διαχειριστής Ουρών (Queue Manager)	31
2 Η ΜΟΝΑΔΑ ΔΙΑΧΕΙΡΙΣΤΗ ΟΥΡΩΝ (QUEUE MANAGER)	32
2.1 Η αρχιτεκτονική του Διαχειριστή Ουρών	33
2.2 Ροές	34
2.3 Ομάδες Ροών (Flow Groups)	36
2.4 Οι εντολές του Διαχειριστή ουρών	37
2.5 Μορφή της Εγγραφής Ροής	39
2.6 Μορφή Εγγραφής ομάδας ροών (Flow Group Record)	40
2.7 Μορφή και ευθυγράμμιση κελιού στη μνήμη SDRAM	42
2.8 Οργάνωση της μνήμης SDRAM	42
2.9 Μηχανές Καταστάσεων (State Machines)	44

2.10	Παράκαμψη λίστας ελευθέρων (Free List Bypassing) [13]	46
2.11	Προανάθεση Χώρου κελιών (Cell Buffer Pre-allocation) [10]	47
2.12	Θέματα χρονισμού	48
2.12.1	Το ρολόι UTOPIA σε σχέση με το ρολόι του Διαχειριστή ουρών (ABRSU).	48
2.12.2	Αποτελέσματα σύνθεσης, συνεισφορά της παράκαμψη λίστας ελευθέρων και της προανάθεσης χώρων κελιού.	49
3	ΣΥΜΠΕΡΑΣΜΑΤΑ ΚΑΙ ΜΕΛΛΟΝΤΙΚΈΣ ΕΠΕΚΤΆΣΕΙΣ	51
	ΠΑΡΑΡΤΗΜΑΤΑ – ΑΓΓΛΙΚΉ ΜΕΤΆΦΡΑΣΗ	53
4	INTRODUCTION	53
4.1	Motivation	53
4.2	This thesis and the DIPOLO Switch	54
4.3	Switch/Router Generations and Queueing Architectures	55
4.3.1	First Generation Switches/Routers	55
4.3.2	Second Generation Switches/Routers	56
4.3.3	Third Generation of Switches/Routers	57
4.4	Queueing Architectures in general	58
4.4.1	Output Queueing	58
4.4.2	Input queueing	59
4.4.3	Variations	60
4.4.4	Per-Flow Queueing Vs Single FIFO Queueing in Output Queueing	61
5	THE DIPOLO ATM SWITCH	63
5.1	The DIPOLO Architecture	63
5.2	The ATM 155Mbps Card	64
5.2.1	The switching device (Transwitch Cubit Pro) [18]	64
5.2.2	The Cell Processing device (Motorola MC92501 Cell Processor) [17]	65
5.2.3	The local CPU (MPC860 by Motorola) [17]	65
5.2.4	The Physical level device (PM5350 S/UNI-155-ULTRA) [20]	65
5.3	The CPU Card	65
5.3.1	Motorola MPC869SAR (PowerQUICC)	66
5.3.2	The CubitPro Device	67
5.3.3	Memories	67
5.4	The ATM Line Card	67
5.4.1	Framer part - PMC-Sierra PM7344 [20]	68
5.4.2	Line interface circuit part - PMC-Sierra PM4314 [20]	68
5.5	The ABR Server Card	68
5.5.1	The Switching device (Transwitch Cubit Pro)	69
5.5.2	The ABR Server Unit (ABRSU) (EPF10K200EBC600-1 FPGA by Altera) [16]	70
5.5.3	The Memory module (256 MB of SDRAM) [19]	70
5.5.4	The local processor (MPC860 by Motorola)	70
6	THE ABR SERVER UNIT (ABRSU)	71

		13
6.1	The ABR Server Architecture	71
6.2	The CPU Interface	72
6.2.1	MPC 860 Interface Block - Queue Manager Block Interface	73
6.2.2	MPC 860 Interface Block - Cell Scheduler Block interconnection	74
6.3	The Cell Demultiplexor	74
6.4	The UTOPIA Interfaces	75
6.4.1	UTOPIA Input Interface	76
6.4.2	UTOPIA Output Interface	77
6.4.3	The Pulse Synchronizer	78
6.5	The Queue Manager	78
6.6	The Cell Scheduler	78
7	THE QUEUE MANAGER IP	80
7.1	The Queue Manager IP Architecture	81
7.2	Functional Implementation	82
7.2.1	Interfaces	82
7.2.1.1	The Interface with the CPU	83
7.2.1.2	The Cell Demux Interface (Incoming cells)	83
7.2.1.3	The Cell Scheduler – Queue Manager Interface	84
7.2.1.4	The SDRAM – Queue Manager Interface	85
7.2.2	Flows	86
7.2.3	Flow Groups	88
7.2.4	Queue Manager Commands	88
7.2.5	EFCI and RM marking [21]	90
7.3	Design Implementation	91
7.3.1	Flow Record Format	91
7.3.2	Flow Group Record Format	92
7.3.3	Cell Format and Alignment	93
7.3.4	SDRAM Memory Organization	94
7.3.5	State Machines	95
7.3.6	The SDRAM Controller	97
7.3.7	Free List Bypassing [13]	98
7.3.8	Cell Buffer Pre-allocation [10]	100
7.4	Timing Issues	100
7.4.1	UTOPIA clock Vs Queue Manager (ABRSU) clock.	100
7.4.2	Worst case of Enqueue and Dequeue	101
7.4.3	Normal case of Enqueue and Dequeue	102
7.4.4	Synthesis Results, Free-list bypassing and Cell Buffer Pre-allocation contribution	104
8	CONCLUSIONS AND FUTURE WORK	105
9	REFERENCES	107

Λίστα σχημάτων

Figure 1-1: Γενική αρχιτεκτονική του μεταγωγέα ΔΙΠΟΛΟ	21
Figure 1-2: Εσωτερικό διάγραμμα της κάρτας Εξυπηρετητή ABR	22
Figure 1-3 :Μηχανισμός Flow Control Κάρτας Εξυπηρετητή ABR	25
Figure 1-4: Εσωτερικό διάγραμμα του ABRSU	27
Figure 1-5: Διάγραμμα της διεπαφής CPU	28
Figure 1-6: Εσωτερικό διάγραμμα του Cell Demultiplexor	30
Figure 2-1: Το εσωτερικό διάγραμμα του διαχειριστή ουρών	33
Figure 2-2: Η δομή των ροών και η αλλαγές σε αυτή μετά από μία εντολή enqueue και μία εντολή dequeue.	36
Figure 2-3: Οι 64 κυκλικές λίστες των ομάδων ροών	37
Figure 2-4: Τα πεδία και η ευθυγράμμιση της εγγραφής ροής	39
Figure 2-5: Οργάνωση της μνήμης ομάδων ροών και των εγγραφών ομάδων ροών.	41
Figure 2-6: Μορφή και ευθυγράμμιση κελιού στη μνήμη SDRAM	42
Figure 2-7: Καταμερισμός και οργάνωση χώρου της μνήμης SDRAM	43
Figure 2-8: Διάγραμμα καταστάσεων των μηχανών καταστάσεων	44
Figure 2-9: Διάγραμμα καταστάσεων της μηχανής καταστάσεων της εντολής Enqueue	46
Figure 2-10: Υλοποίηση της παράκαμψης λίστας ελευθέρων στον διαχειριστή ουρών	47
Figure 4-1: First Generation Switch Routers	56
Figure 4-2: Second Generation Switch/Router	57
Figure 4-3: Left: Third generation Switch/Router, Top-Right: A crossbar, Bottom-Right: An 8x8 Banyan Fabric made of small 2x2 Switch blocks.	58
Figure 4-4: Left: Output Queueing with a Switching Fabric and multiple buffers Right: Input Queueing with a Switching Fabric.	59
Figure 4-5: Left : Head of Line Blocking Right: Advanced Input Queueing	60
Figure 4-6: Left: Internal Speed Up Switch, Right: Crosspoint Queueing Switch	60
Figure 4-7: Single FIFO queueing and two threshold congestion detection approach	61
Figure 4-8: Per-Flow queueing and two-threshold detection approach	62
Figure 5-1: General Architecture of DIPOLO ATM Switch	63
Figure 5-2: The ATM 155 Card block diagram	64
Figure 5-3: The CPU Card block diagram	66
Figure 5-4: The VDSL Line Card Block Diagram	67
Figure 5-5: ABR Server Card block diagram	69
Figure 6-1: The ABRSU internal block diagram	72
Figure 6-2: The CPU Interface sub-block diagram	73
Figure 6-3: The Cell Demultiplexor sub-block diagram	75
Figure 6-4: UTOPIA Input Interface block diagram	76
Figure 6-5: UTOPIA Output Interface sub-block diagram	77
Figure 6-6: The Pulse Synchronizer	78
Figure 7-1: The Queue Manager IP sub-block diagram	81
Figure 7-2: The Flow Structure and changes on it after an enqueue and a dequeue operation.	87
Figure 7-3: The 64 Flow Group cyclic lists	88
Figure 7-4: Flow Record Fields and alignment	91
Figure 7-5: Flow Group memory organization and Flow Group records.	93
Figure 7-6: Cell Format and alignment	94
Figure 7-7: SDRAM Memory-space division and organization	94
Figure 7-8: State Machine Top Level Diagram	95
Figure 7-9: Enqueue command FSM bubble diagram	97
Figure 7-10 : State Machine Diagram of the SDRAM controller	98
Figure 7-11: Free List Bypassing implementation in the Queue Manager IP	99
Figure 7-12: Worst case Enqueue-Dequeue timing Diagram	102
Figure 7-13: Normal case Enqueue-Dequeue timing diagram	103

Λίστα πινάκων

Table 2-1: Πίνακας όλων των εντολών του Διαχειριστή Ουρών	38
Table 2-2 : Πεδιά της εγγραφής ροής και περιγραφή τους	39
Table 2-3: Περιγραφή των πεδίων των εγγραφών ομάδων ροών	41
Table 2-4: Προτεραιότητες των εντολών του Διαχειριστή ουρών	45
Table 7-1: The Interface with the CPU	83
Table 7-2 : The Cell Demux Interface (Incoming cells)	83
Table 7-3: The Cell Scheduler – Queue Manager Interface	84
Table 7-4: The SDRAM – Queue Manager Interface	86
Table 7-5: Table of all the Queue Manager commands	89
Table 7-6 : Flow Record Field bits and description	91
Table 7-7: Flow Record Field description	93
Table 7-8: Queue Manager Command Priorities	96

1 Εισαγωγή

1.1 Κίνητρα

Η εισαγωγή των εφαρμογών με μεγάλες απαιτήσεις εύρους διαμεταγωγής στις εταιρικές και πελατιακές επικοινωνίες είναι μια από τις πιο σταθερές τάσεις στον δικτυακό χώρο. Οι εφαρμογές πολυμέσων που ακολουθούν το Νόμο του Moore παράγουν υπερβολικές ποσότητες δεδομένων εικόνας και ήχου προς αναμετάδοση πάνω από το διαδίκτυο, ενώ η ραγδαία αυξανόμενη δικτύωση επιχείρησης με επιχείρηση (Business to Business networking) συνεισφέρει μικρά αλλά συχνά πακέτα δεδομένων στο δίκτυο ανάμεσα στις κεντρικές έδρες τους. Αυτές οι δικτυακές απαιτήσεις είναι βαρύ φορτίο για την υποδομή των σύγχρονων δικτύων που βασίζεται κυρίως στο πρωτόκολλο IP και στα καλώδια χαλκού. Ενώ τα δεύτερα είναι το κύριο αντικείμενο του προβλήματος τελευταίου μιλίου (last-mile problem) και απαιτεί σταδιακή αλλά γιγαντιαία επένδυση κεφαλαίων σε παγκόσμιο επίπεδο, για την εισαγωγή οπτικών ινών, το πρωτόκολλο IP είναι ανίκανο να διαφοροποιήσει τις διαφορετικές υπηρεσίες που απαιτούνται από τους δικτυακούς χρήστες. Το πρωτόκολλο ATM θα πρέπει σταδιακά να αντικαταστήσει ή να συγχωνευθεί με τις εφαρμογές τύπου IP στην πορεία για υψηλής ταχύτητας, οπτικά δίκτυα.

Τα δίκτυα ATM προσφέρουν εγγυήσεις ποιότητας υπηρεσιών (Quality of Service (QoS) guarantees) στα σύγχρονα δίκτυα με το να διαφοροποιούν την δικτυακή κίνηση σε διαφορετικά είδη με βάση τις απαιτήσεις των εφαρμογών που την γεννούν και παρέχει ειδικό χειρισμό και χρέωση για κάθε είδος.

Αυτές οι απαιτήσεις είναι:

- **Εύρος Διαμεταγωγής (Bandwidth)** – Ο ρυθμός με τον οποίο το δίκτυο πρέπει να μεταφέρει την κίνηση μιας εφαρμογής.
- **Καθυστέρηση (Delay)** – Η καθυστέρηση που μια εφαρμογή μπορεί να υποστεί στην παράδοση των δεδομένων της.
- **Απόκλιση (Jitter)** – Η διαφοροποίηση στην καθυστέρηση.
- **Απώλειες (Loss)** – Το ποσοστό της αποδεκτής απώλειας δεδομένων.

Κάθε μία από την οικογένεια των ειδών κίνησης που υποστηρίζουν τα δίκτυα ATM δίνει πολύ έμφαση σε κάποιες από τις παραπάνω απαιτήσεις και λίγη στις υπόλοιπες. Τα είδη της κίνησης ATM είναι:

- **CBR (Constant Bit Rate)** – Αυτό το είδος της κίνησης χρησιμοποιείται για την προσομοίωση της μεταγωγής κλειστού κυκλώματος (Circuit switching). Το εύρος μεταγωγής είναι σταθερό στο χρόνο. Οι εφαρμογές τύπου CBR είναι ευαίσθητες στην απόκλιση (Jitter) αλλά όχι τόσο πολύ στην απώλεια δεδομένων. Παραδείγματα τέτοιων εφαρμογών είναι η κίνηση τηλεφωνείας, η βιντεο-συνδιάσκεψη και η τηλεόραση υψηλής ευκρίνειας.
- **VBR-NRT (Variable Bit Rate – Non Real Time)** – Αυτό το είδος κίνησης επιτρέπει στους χρήστες να στέλνουν κίνηση με ρυθμό που κυμαίνεται στο χρόνο ανάλογο με την διαθεσιμότητα της πληροφορίας του χρήστη. Χρησιμοποιείται η στατιστική πολύπλεξη για να κάνει την βέλτιστη χρήση των δικτυακών πόρων. Το πολυμεσικό ηλεκτρονικό ταχυδρομείο είναι ένα παράδειγμα εφαρμογής τύπου VBR-NRT.

- VBR-RT (Variable Bit Rate – Real Time) – Αυτό το είδος της κίνησης είναι παρόμοιο με το VBR-NRT αλλά είναι σχεδιασμένο για εφαρμογές που είναι ευαίσθητες στη διαφοροποίηση της απόκλισης της καθυστέρησης κελιών (Cell delay variation). Παραδείγματα εφαρμογών τύπου VBR-RT είναι ήχος με ενεργοποίηση δραστηριότητας (Voice with speech activity detection) και αλληλεπιδραστικό συμπιεσμένο βίντεο.
- ABR (Available Bit Rate) – Η κίνηση αυτού του είδους παρέχει έλεγχο ροής και στοχεύει στην κίνηση δεδομένων όπως είναι η μεταφορά αρχείων και το e-mail. Παρόλο που οι προδιαγραφές της δεν απαιτούν η καθυστέρηση της μετάδοσης των κελιών και ο λόγος απώλειας κελιών να εγκυάται ή να περιορίζεται, είναι επιθυμητό να είναι όσο περιορισμένα γίνεται. Ανάλογα με τα επίπεδα της συμφόρησης στο δίκτυο, η πηγή εξαναγκάζεται να περιορίσει το ρυθμό αποστολής των δεδομένων της. Οι χρήστες δικαιούνται να δηλώσουν τον ελάχιστο δυνατό ρυθμό αποστολής, που παρέχεται στη σύνδεση από το δίκτυο.
- UBR (Unspecified Bit Rate) – Αυτό το είδος παρέχεται για όλες τις άλλες περιπτώσεις εφαρμογών που δεν έχουν καμία απαίτηση και χρησιμοποιείται ευρύτατα σήμερα για μετάδοση TCP/IP.

Τα τελευταία δύο είδη κίνησης είναι τα πιο ομοιογενή, μιας και τείνουν να χρησιμοποιούν τα περισεύματα των δικτυακών πόρων που έχουν δεσμευθεί από τις υπόλοιπες, οπότε αυτό είναι εφικτό. Παρόλα αυτά, μεταφέρουν το βάρος των πιο παραδοσιακών δικτυακών εφαρμογών όπως Ηλεκτρονικό Ταχυδρομείο, FTP, HTML που είναι η βάση της σύγχρονης μορφής δικτύων. Αν και θέτουν κάποιες απαιτήσεις στο ρυθμό διαμεταγωγής, παραμερίζονται σε περιόδους δικτυακής συμφόρησης. Αυτό μεταφράζεται σε εκτεταμένη αποθήκευση αυτών των ειδών στους δικτυακούς κόμβους που βρίσκονται, όταν οι τελευταίοι πάσχουν από συμφόρηση. Υπάρχει μία σειρά από προτεινόμενες αρχιτεκτονικές σε όλη την εξέλιξη των μεταγωγέων, που αντιμετωπίζει το πρόβλημα της αποδοτικής αποθήκευσης της κίνησης σε ουρές. Οι μεγάλοι αποθηκευτικοί χώροι από μόνοι τους δεν μπορούν να υποστηρίξουν τη ποικιλία των ειδών κίνησης και τα προβλήματα των καθυστερήσεων και του προβλήματος HoL (Head of Line) που εμφανίζονται. Προκύπτει ότι η παροχή ξεχωριστής αποθήκευσης για κάθε μία ροή κίνησης (per-flow queueing) εκμεταλλεύεται τον αχρησιμοποίητο χώρο των ανενεργών ροών [5]. Οι μεταγωγείς με τέτοια χαρακτηριστικά μπορούν να υποστηρίξουν κίνηση δικτύων υψηλών ταχυτήτων χρησιμοποιώντας μονάδες δυναμικής μνήμης RAM χαμηλού κόστους, ρίχνοντας το κόστος της μνήμης στο ελάχιστο επίπεδο αυτής ενός προσωπικού υπολογιστή.

Η δυναμική Μνήμη PAM σε όλες τις μορφές της έχει ένα μεγάλο εύρος από εφαρμογές στο σύγχρονο χώρο των υπολογιστών και ειδικά οι σύγχρονες δυναμικές μνήμες (SDRAM) που χρησιμοποιούνται σχεδόν σε κάθε υπολογιστική συσκευή που χρειάζεται μεγάλο και χαμηλού κόστους αποθηκευτικό χώρο, με μεγάλο ρυθμό προσπέλασης δεδομένων. Οι καλύτερες επιδόσεις επιτυγχώνονται όταν μεγάλα πακέτα δεδομένων προσπελαύνονται, κάτι που κάνει αυτές τις μνήμες ιδανικές για δικτυακές συσκευές που βρίσκονται εδώ και καιρό στη εποχή της μεταγωγής πακέτων (Packet Switching). Η ευρύτατη εφαρμογή τους, η απλοϊκή διεπαφή τους και η βιομηχανική τυποποίηση τους ρίχνουν το κόστος τους σημαντικά, τροφοδοτώντας της πρόσθεση επιπλέον δικτυακών πορτών στις δικτυακές συσκευές [1, κεφ 9], [12].

1.2 Περιεχόμενα αυτής της εργασίας στο πλαίσιο του προγράμματος ΔΙΠΟΛΟ.

Σε αυτή την εργασία περιγράφουμε την αρχιτεκτονική ενός Διαχειριστή Ουρών που υποστηρίζει τα χαρακτηριστικά της ανά-ουρά δυναμικής αποθήκευσης σε ουρά κίνησης τύπου ABR δικτών ATM. Αυτός ο διαχειριστής ουρών σχεδιάστηκε για τις ανάγκες του μεταγωγέα ATM ΔΙΠΟΛΟ, ικανότητας διαμεταγωγής 1 Gbps.

Χρησιμοποιείται για την παροχή αποθήκευσης σε 64 χιλιάδες ροές μιας κεντρικοποιημένης κάρτας παροχής ABR. Η κάρτα χρησιμοποιεί μια μεγάλη FPGA για την στέγαση του Διαχειριστή Ουρών και μια μοναδική μονάδα SDRAM DIMM για την αποθήκευση των κελιών και των δεικτών των ουρών. Ακόμη υπάρχει μια διεπαφή επεξεργαστή που προγραμματίζει τον Διαχειριστή ουρών με παραμέτρους ροής τύπου ABR ώστε ο τελευταίος να μπορεί να υποστηρίξει χαρακτηριστικά ελέγχου ροής τύπου ABR (ABR Flow Control), όπως μαρκάρισμα RM και EFCI. Ακόμη ασχολούμαστε με τεχνικές που χρησιμοποιήθηκαν στον διαχειριστή για να αυξήσουν την χρήση της μνήμης, όπως την προανάθεση μνήμης (Buffer Preallocation) και την παράκαμψη ελεύθερης λίστας (Free List Bypassing). Ο Διαχειριστής ουρών υποστηρίζει διαφοροποίηση κίνησης, με το να οργανώνει τις ροές σε ομάδες ροών με βάση τις ανάγκες εξυπηρέτησής τους και την πόρτα εξόδου τους, κάτι που κάνει εφικτή την υποστήριξη κίνησης τύπου CBR και VBR, όταν υπάρχει ειδική φροντίδα στο υλικό χρονοπρογραμματισμού (scheduler hardware). Ειδική φροντίδα έχει ληφθεί ώστε οι διεπαφές του Διαχειριστή με τις άλλες μονάδες όπως ο χρονοπρογραμματιστής (Scheduler) και ο Επεξεργαστής, να είναι όσο το δυνατόν πιο απλές και αποδοτικές. Η μοναδική μονάδα μνήμης SDRAM δεν επιτρέπει παράλληλες προσπελάσεις. Αντ' αυτού οι εντολές προσθήκης κελιού σε ουρά (Enqueue command) και αφαίρεσης κελιού από την κορυφή μιας ουράς (Dequeue) υλοποιήθηκαν ενιαία, αντί σε μικρές απλές εντολές, ώστε να βελτιωθεί η χρησιμοποίηση της μνήμης. Πετύχαμε συχνότητα ρολογιού 35 MHz που μεταφράζεται σε 800 Mbps μέγιστης εισερχόμενης και εξερχόμενης διαμεταγωγής. Σε υλοποιήσεις τύπου ASIC όπου οι συχνότητες ρολογιού της τάξης των 133 MHz, η διαμεταγωγή που πετύχαμε μπορεί να αυξηθεί κάνοντας την αρχιτεκτονική μας κατάλληλη για χρήση ως μέρος ενός δικτυακού ολοκληρωμένου [7], [10]. Ο διαχειριστής ουρών χρησιμοποιεί 2500 λογικά στοιχεία της FPGA και 2000 bit εσωτερικής SRAM, για 64000 χιλιάδες ουρές.

Ο μεταγωγέας ΔΙΠΟΛΟ σχεδιάστηκε από κοινού από το Πανεπιστήμιο Κρήτης, το Εθνικό Μετσόβιο Πολυτεχνείο, το Ίδρυμα Τεχνολογίας και Έρευνας της Κρήτης και την Ιντρακομ στα πλαίσια του έργου Ε.Π.Ε.Τ «Δίκτυα Πρόσβασης Ολοκληρωμένων Υπηρεσιών (ΔΙΠΟΛΟ)». Σκοπός της κοινής δραστηριότητας ήταν η σχεδίαση και η κατασκευή ενός μεταγωγέα ATM με ικανότητα διαμεταγωγής 1 Gbit το δευτερόλεπτο σε οικιακούς χρήστες γραμμών VDSL. Απαραίτητη ήταν η υποστήριξη από τον μεταγωγέα, κίνησης τύπων CBR, VBR και ABR. Ακόμη ελήφθησαν υπόψιν παράμετροι όπως το χαμηλό κόστος του συστήματος, η χρήση εμπορικών ολοκληρωμένων και ο βέλτιστος καταμερισμός των απαραίτητων εργασιών ανάμεσα στους συμμετέχοντες οργανισμούς.

Στα πλαίσια της σχεδίασης του μεταγωγέα ΔΙΠΟΛΟ, ο συγγραφέας είχε την ευκαιρία να ασχοληθεί με την φυσική σχεδίαση της κάρτας Εξυπηρετητή ABR. Η κάρτα αυτή περιγράφεται στην παράγραφο 1.4. Πιο συγκεκριμένα ασχολήθηκε με τον καθορισμό της οργάνωσής της κάρτας, την περιγραφή του σε επίπεδο CONCEPT. Παράλληλα

καθόρισε την εσωτερική οργάνωση της μονάδας ABRSU τύπου FPGA της κάρτας. Η μονάδα περιγράφεται στην παράγραφο 1.5. Σχεδίασε τις λειτουργικές μονάδες Cell Demux (Δες παράγραφο 1.5.4) τον Διαχειριστή ουρών, κύριο θέμα αυτής της εργασίας, που περιγράφεται εκτενώς στο κεφάλαιο 2, την διεπαφή όλων των μονάδων με τον μικροελεγκτή της κάρτας (Δες παράγραφο 1.5.1-3), καθώς και την διεπαφή με την μονάδα μνήμης της Κάρτας. Ακόμη ήταν υπεύθυνος για την ολική σύνθεση της μονάδας και την σχεδίαση και εκτέλεση δοκιμαστικών πειραμάτων και επιδείξεων όλης της κάρτας για την τεκμηρίωση των δραστηριοτήτων.

1.3 Η αρχιτεκτονική του μεταγωγέα ΔΙΠΟΛΟ

Η αρχιτεκτονική του μεταγωγέα ATM ΔΙΠΟΛΟ παρουσιάζεται στο σχήμα 1-1. Όλο το σύστημα αποτελείται από μια σειρά από κάρτες γραμμής τύπου VDSL που συνδέουν τους χρήστες με το μεταγωγέα, μία ή περισσότερες κάρτες ATM 155 για τη διασύνδεση του μεταγωγέα με το κεντρικό δίκτυο ATM μέσω γραμμής OC-3/STM-1, μια κάρτα Επεξεργαστή, μια κάρτα Διαχειριστή ABR (σε αυτή υλοποιείται και ο διαχειριστής ουρών που είναι το κύριο θέμα αυτής της εργασίας) και δύο δίαυλοι CellBus.

Η ακριβής λειτουργικότητα κάθε κάρτας δίνεται παρακάτω:

- Κοινός δίαυλος κελιών: Σε επίπεδο συστήματος, η μεταφορά των δεδομένων όσο και της πληροφορίας σηματοδοσίας και διαχείρισης γίνεται μέσω συνδέσεων ATM. Η μεταγωγή των δεδομένων αυτών των συνδέσεων πραγματοποιείται μέσω της χρήσης ενός κοινού δίαυλου κελιών (Cellbus). Για λόγους ασφαλείας η αρχιτεκτονική χρησιμοποιεί δύο δίαυλους, έναν πρωτεύοντα και ένα δευτερεύοντα που θα χρησιμοποιείται μόνο σε περίπτωση αστοχίας του πρωτεύοντος.
- Κάρτα ATM155: Εκτελεί όλες τις λειτουργίες που εμπεριέχονται στο τέλος του φυσικού επιπέδου και σε ολόκληρο το επίπεδο ATM όπως τερματισμός του φυσικού επιπέδου, δρομολόγηση ATM, σηματοδοσία Q.2931, λειτουργίες OA&M, κ.τ.λ. Η σύνδεση πολλών τέτοιων καρτών σε ένα backplane οδηγεί στην υλοποίηση ενός μεταγωγέα.
- Εξυπηρετητής ABR κίνησης: Εκτελεί όλες τις απαραίτητες λειτουργίες για την αποδοχή, αποθήκευση και προώθηση των κελιών που μεταφέρουν την κίνηση ABR την οποία δεν μπορεί να χειριστεί η κάρτα ATM γιατί απαιτεί μεγάλο αποθηκευτικό χώρο. Η βασικότερη από αυτές τις λειτουργίες είναι η διαχείριση μεγάλου αριθμού ουρών πολλαπλών επιπέδων.
- Κεντρικός επεξεργαστής: Έχει την συνολική επίβλεψη/διαχείριση του συστήματος και θα εκτελεί λειτουργίες υψηλού επιπέδου όπως το Call Admission Control (CAC).
- Κάρτα VDSL modem: Υλοποιεί την λειτουργικότητα των modem τύπου VDSL για την σύνδεση με τους χρήστες μέσω των τηλεφωνικών καλωδίων. Η οργάνωση της κάρτας αυτής είναι έξω από τον σκοπό της αναφοράς αυτής.

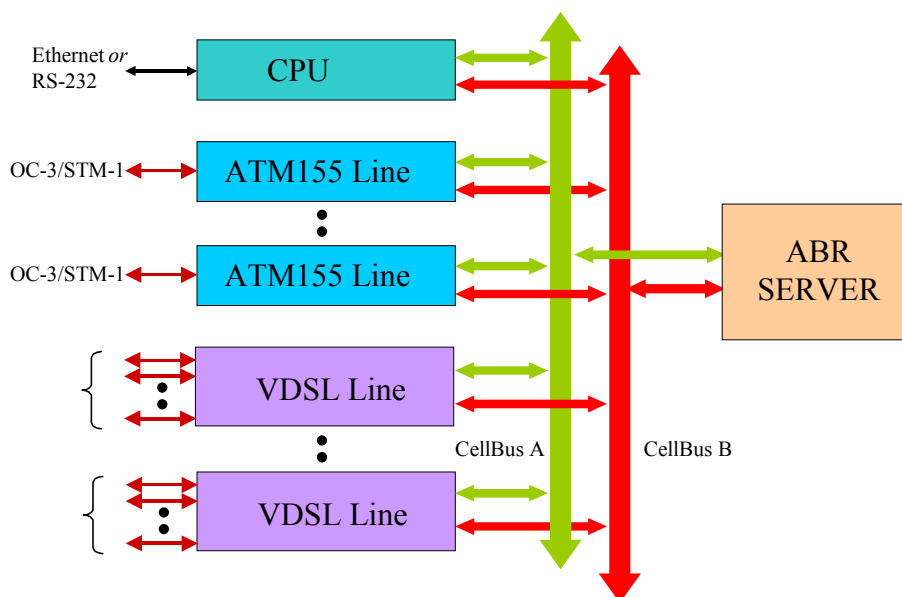


Figure 1-1: Γενική αρχιτεκτονική του μεταγωγέα ΔΙΠΟΛΟ

Στην οργάνωση αυτή, η κεντροποιημένη εξυπηρέτηση της κίνησης ABR επιτρέπει την συγκέντρωση της μνήμης προσωρινής αποθήκευσης των δεδομένων σε μία κάρτα, με συνέπεια την καλύτερη χρήση της μνήμης αυτής. Ακόμα επιτρέπει την δημιουργία συστημάτων με ή χωρίς εξυπηρετητή ABR, και αυξάνει την ευελιξία του συστήματος στην αποδοτική (από πλευράς κόστους) κάλυψη μεταβαλλόμενων αναγκών.

1.4 Αρχιτεκτονική του εξυπηρετητή ABR κίνησης

Η ABR υπηρεσία των δικτύων ATM έρχεται να καλύψει τις ανάγκες των παραδοσιακών δικτύων LAN. Τα υπολογιστικά συστήματα σε ένα LAN θέλουν να στείλουν τα δεδομένα τους, τη στιγμή που αυτά γενιούνται και με ταχύτητα, αν αυτό είναι δυνατό, αυτή της γραμμής αλλά χωρίς συνωστισμό (congestion) που προκαλεί απώλεια cells. Αυτό ισχύει γιατί τα δεδομένα των Ηλεκτρονικών Υπολογιστών είναι ευαίσθητα στις απώλειες, μιας και οι αναμεταδώσεις μπορεί να μειώσουν δραστικά την απόδοση όλου του δικτύου.

Αντί λοιπόν να δεσμεύονται πόροι για την εκρηκτική κίνηση (bursty traffic) των LAN δικτύων, αυτή εξυπηρετείται από την υπηρεσία ABR που χρησιμοποιεί την διαθέσιμη κίνηση που δεν χρησιμοποιούν οι υπόλοιπες, ευαίσθητες σε χρονισμό, υπηρεσίες του ATM δικτύου. Προκειμένου όμως να δοθεί στους χρήστες η ικανότητα να στέλνουν ότι θέλουν όταν το θέλουν, με μόνη απαίτηση να μην υπάρχουν απώλειες, η εξυπηρέτηση της κίνησης ABR πρέπει να γίνεται από ένα μεταγωγέα που διαθέτει μεγάλη αποθηκευτική ικανότητα, ώστε να μπορούν να αποθηκευτούν για κάποιο χρονικό διάστημα τα cells που δεν μπορούν να μεταδοθούν λόγω μη διαθέσιμης κίνησης.

Γι' αυτό, στο μεταγωγέα που σχεδιάζουμε, χρησιμοποιούμε ένα κεντροποιημένο εξυπηρετητή κίνησης ABR, ώστε να μπορούμε να εκμεταλευτούμε ικανοποιητικά τους πόρους μνήμης με χαμηλό κόστος. Πιο συγκεκριμένα η κάρτα εξυπηρέτησης ABR κίνησης διαθέτει ικανότητες αποθήκευσης της συνολικής κίνησης ABR που δέχεται ο μεταγωγέας. Με αυτό τον τρόπο δεν παρένοχλείται η εξυπηρέτηση των άλλων ειδών κίνησης και η μνήμη της κάρτας χρησιμοποιείται αποκλειστικά για

κίνηση ABR, που προωθείται όταν ο διάυλος του μεταγωγέα είναι ελεύθερος (όταν δηλαδή δεν υπάρχει άλλου είδους κίνησης) απλοποιώντας έτσι την διαχείριση της κίνησης στις κάρτες γραμμής. Έτσι η κεντρικοποιημένη εξυπηρέτηση ABR μειώνει το κόστος μνήμης μιας και δεν τοποθετούμε μνήμη σε κάθε κάρτα γραμμής.

Αν και η κεντρικοποιημένη εξυπηρέτηση ABR δίνει μοναδικό σημείο αποτυχίας, το πρόβλημα μπορεί να λυθεί με την χρήση δεύτερης κάρτας που θα παρέχει ανοχή σε σφάλματα και μοίρασμα φόρτου.

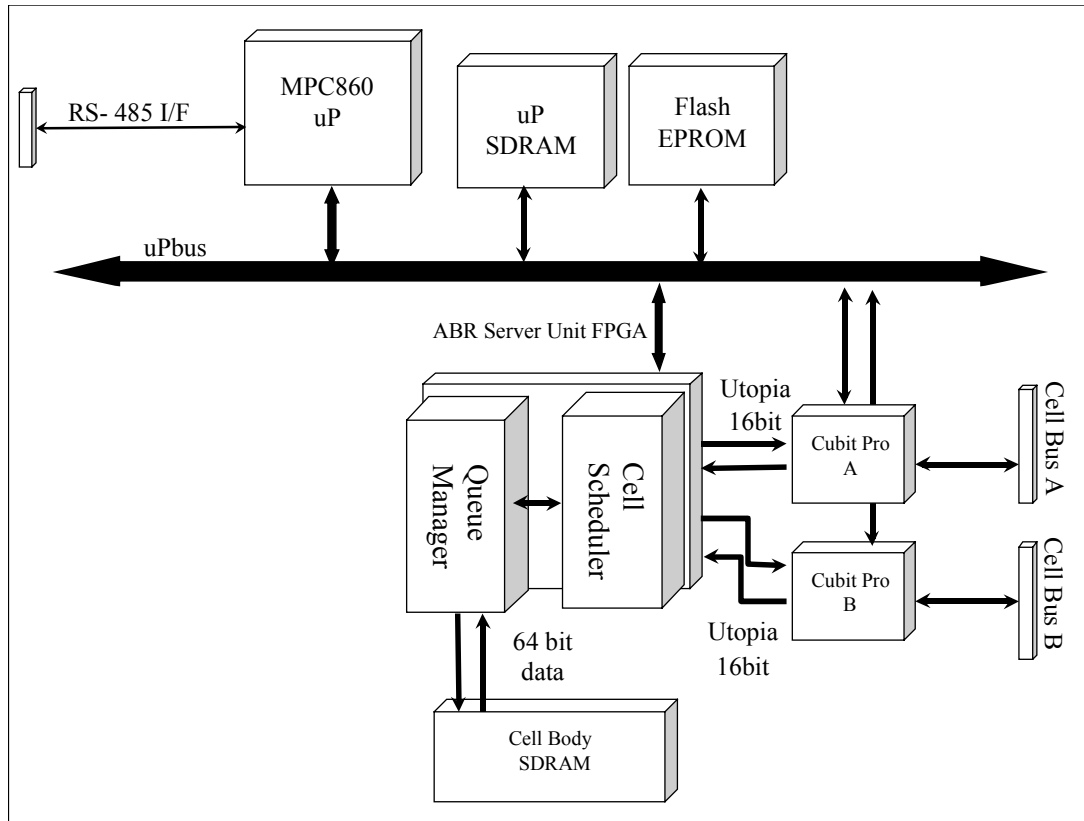


Figure 1-2: Εσωτερικό διαγραμμα της κάρτας Εξυπηρετητή ABR

Η κάρτα εξυπηρέτησης ABR αποτελείται από τα ακόλουθα στοιχεία:

1. Τη Συσκευή Μεταγωγής (Transwitch Cubit Pro)
2. Τη Μονάδα Εξυπηρέτησης ABR (The ABR Server Unit (ABRSU))
3. Τον Επεξεργαστή MPC860.
4. The Memory module (256 MB of SDRAM)

Η γενική αρχιτεκτονική και η οργάνωση της κάρτας φαίνονται στο σχήμα 1.2.

1.4.1 Το στοιχείο μεταγωγής

Το στοιχείο μεταγωγής (Transwitch Cubit Pro) συνδέεται με την Μονάδα Εξυπηρέτησης ABR (ABR Server Unit), μέσω μιας σύνδεσης UTOPIA και με το Backplane μέσω των 37 γραμμών του CellBus. Ακόμη συνδέεται με τον επεξεργαστή MPC860, μέσω μιας διασύνδεσης επεξεργαστή (με επιλεγμένη την κατάσταση Motorola). Κάθε κάρτα εξυπηρέτησης ABR διαθέτει 2 Cubit Pro, με την κάθε μια να

είναι συνδεδεμένη σε ένα από τα δύο CellBus. Μόνο ένα CubitPro είναι σε λειτουργία μια δεδομένη στιγμή ενώ το άλλο είναι σε αναμονή για την περίπτωση αστοχίας του πρωτεύοντος CellBus. Σε αυτή τη περίπτωση το εν αναμονή CubitPro τίθεται σε λειτουργία ενώ αυτό του πρωτεύοντος CellBus απενεργοποιείται. Το Cubit-Pro δέχεται Cells από το CellBus και τα προωθεί στη Μονάδα Εξυπηρέτησης ABR μέσω μιας ουράς FIFO 123 cells. Στην αντίθετη κατεύθυνση αποδέχεται cells από τη στη Μονάδα Εξυπηρέτησης ABR και τα προωθεί στο CellBus μέσω μιας ουράς FIFO 4 cells. Η συσκευή CubitPro παρέχει μαρκάρισμα ένδειξης Συνοστισμού (Congestion indicator marking) και πέταγμα cells (cell discarding) κάτω από ορισμένες συνθήκες.

1.4.2 Εσωτερικός επεξεργαστής ελέγχου

Ο επεξεργαστής MPC860 βασίζεται στην αρχιτεκτονική PowerQUICC της Motorola. Είναι υπεύθυνος για την ορθή λειτουργία, αρχικοποίηση και ανίχνευση λαθών της κάρτας εξυπηρέτησης ABR. Συνδέεται με διάφορες μονάδες της κάρτας μέσω ενός διαύλου uPBus. Χρησιμοποιεί την διασύνδεση μικροεπεξεργαστή των συσκευών CubitPro για να ανταλλάξει cells διαχείρισης (Management Cells) με τον μικροεπεξεργαστή που βρίσκεται στη κάρτα CPU. Το CubitPro παρέχει (μέσω polling ή interrupts), επίσης, πληροφορίες για την κίνηση και στατιστικά, χρήσιμα για την διαχείριση του όλου συστήματος. Ακόμη σε περίπτωση ελαττωματικής λειτουργίας του ενεργού CellBus, ο MPC860 είναι υπεύθυνος για την απενεργοποίησή του και την ενεργοποίησή του αναπληρωματικού. Ακόμη, ο MPC860 διασυνδέεται με την ABRSU και λαμβάνει QIDs για την αρχικοποίηση μιας κλήσης σύνδεσης, ενώ σε περίπτωση τερματισμού μιας σύνδεσης επιστρέφει το αντίστοιχο QID [22].

Ο επεξεργαστής MPC860 Χρησιμοποιεί μια μονάδα Flash EPROM που αποθηκεύει δεδομένα για την αρχικοποίηση του συστήματος και μια SDRAM μνήμη για αποθήκευση δεδομένων και cells. Λόγω των ικανοτήτων ελεγκτή μνήμης (on-chip memory controller), ο MPC860 μπορεί να συνδεθεί με τα παραπάνω χωρίς επιπλέον λογική.

Τέλος, για τον έλεγχο του συστήματος, συντηρεί μια σύνδεση RS-232/423, που ελέγχεται και αυτή μέσω ενός σειριακού ελεγκτή πάνω στο ολοκληρωμένο του MPC860.

1.4.3 Η μονάδα εξυπηρέτης κίνησης ABR

Η μονάδα εξυπηρετητή κίνησης ABR (ABR Server Unit – ABRSU) αποτελείται από 2 στοιχεία, τον προγραμματιστή Cell (Cell Scheduler) και τον διαχειριστή ουρών (Queue Manager). Συνδέεται με τις συσκευές CubitPro μέσω μιας αμφίδρομης διασύνδεσης UTOPIA πλάτους 16 bit data, με τον επεξεργαστή MPC860 και με τον Cell Processor μέσω μιας αμφίδρομης διασύνδεσης UTOPIA πλάτους 8 bit data. Μια μονάδα SDRAM είναι συνδεδεμένη με την ABRSU με σκοπό την αποθήκευση Cells καθώς και δομών δεδομένων που υλοποιούν αυτές τις ουρές. Η λειτουργία της ABRSU συνίσταται στη αποθήκευση των cells τύπου ABR με βάση ένα πεδίο 16 bit (εισερχόμενη κατεύθυνση) βρίσκεται στον header κάθε cell και ορίζει μοναδικά την εικονική σύνδεση που αυτό ανήκει, και στον προγραμματισμό μετάδοσης των cells (εξερχόμενη διεύθυνση). Το σύνολο της λογικής της μονάδας εξυπηρέτησης κίνησης ABR υλοποιείται από μια FPGA συνδεδεμένη με μνήμη SDRAM

για την αποθήκευση των δεδομένων και των εσωτερικών δομών (head-tail pointers για τις ουρές, κ.α.)

Η μονάδα εξυπηρέτησης κίνησης ABR υποστηρίζει 64K ουρές (αφού το πεδίο στον header είναι 16 bits), ώστε να αποφευχθούν προβλήματα «head-of-line blocking», και υποστηρίζει αρκετό αποθηκευτικό χώρο για τα δεδομένα των cells ώστε να μπορεί να απορροφά διακυμάνσεις στην διαθεσιμότητα bandwidth.

Αν και το σύστημα που θα υλοποιήσουμε σε αυτό το έργο θα εξυπηρετεί σχετικά μικρό αριθμό χρηστών (το πολύ λίγες δεκάδες) και αυτή την περίπτωση η υποστήριξη ακόμα και λίγων εκατοντάδων ουρών θα ήταν αρκετή για την υλοποίηση πολλαπλών συνδέσεων με κίνηση ABR ανά χρήστη, τελικά, η μονάδα εξυπηρέτησης κίνησης ABR θα υποστηρίζει 64 χιλιάδες ουρές, έτσι ώστε να μπορεί να χρησιμοποιηθεί χωρίς αλλαγές και σε μεγαλύτερα συστήματα.

Από πλευράς χώρου για την αποθήκευση των δεδομένων, δεδομένης της χρήσης συνδέσεων με ταχύτητα 155 Mbps, η χρήση μνήμης μεγέθους λίγων δεκάδων Mbytes αρκεί για την αποθήκευση του συνόλου της κίνησης για μερικά δευτερόλεπτα. Δεδομένου ότι η μνήμη θα εξυπηρετεί μόνο κίνηση τύπου ABR ακόμα και λίγα Mbytes θα είναι ικανά να ανταποκριθούν σε σημαντικές διακυμάνσεις της διαθεσιμότητας του bandwidth. Η δε χρήση μνημών τύπου SDRAM επιβάλλει ένα ελάχιστο μέγεθος 256 Mbytes, το οποίο θα είναι υπεραρκετό για τις ανάγκες της κίνησης ABR.

1.4.4 Η μονάδα μνήμης SDRAM (256 MB) [19]

Μια μονάδα μνήμης τύπου DIMM SDRAM χρησιμοποιήθηκε ως εξωτερικός αποθηκευτικός χώρος της Μονάδας Εξυπηρετητή ABR (ABRSU), όπου αποθηκεύονται κελιά σε ουρές καθώς και άλλα δεδομένα ουρών και πληροφορίες ελέγχου ροής. Αν και λιγότερα από 256 MB είναι αρκετά για την κάρτα μας, τα DIMM μεγέθους 256 MB είναι αρκετά συνήθη στην αγορά.

1.4.5 Μηχανισμός Flow Control Κάρτας Εξυπηρετητή ABR [23]

Όπως αναφέρεται και στο 5.2.1, το CellBus και τα ολοκληρωμένα CubitPro έχουν τη δυνατότητα να μεταφέρουν πληροφορία Flow Control αναμεταξύ των καρτών που χρησιμοποιούν το CellBus για μεταγωγή cells. Αυτό γίνεται μέσω ειδικών σημάτων του CellBus που θέτουν και διαβάζουν τα CubitPro, ώστε με τη σειρά τους να ειδοποιήσουν τα υπόλοιπα ολοκληρωμένα της κάρτας τους για την ύπαρξη congestion σε κάποια από τις κάρτες ή τη μη ορθή λήψη ενός cell από μια κάρτα προορισμού.

Τον παραπάνω μηχανισμό Flow Control χρησιμοποιεί και η κάρτα Εξυπηρετητή ABR για την ορθή μεταγωγή cells στις κάρτες προορισμού καθώς και για να πληροφορηθεί για ποιές κάρτες δεν διαθέτουν χώρο για λήψη ενός cell τύπου ABR στην αντίστοιχη ουρά του CubitPro τους. Έτσι, όσες συνδέσεις που εξυπηρετεί η κάρτα στέλνουν cells στις παραπάνω κάρτες δεν περιλαμβάνονται στις πιθανές για μετάδοση συνδέσεις από τον Cell Scheduler.

Ο μηχανισμός φαίνεται στο σχήμα 5. Όπως φαίνεται στο σχήμα το CubitPro της κάρτας Εξυπηρετητή ABR μεταδίδει Flow Control πληροφορία μέσω τριών σημάτων:

- CONGOUT : Όταν το σήμα αυτό είναι ενεργό σημαίνει ότι στη κάρτα στην οποία στάλθηκε το τελευταίο cell τύπου ABR υπάρχει congestion. Έτσι ο Cell Scheduler του ABRSU θα πρέπει να σταματήσει να στέλνει, για κάποιο χρονικό διάστημα, πάνω από το CellBus, cells που προορίζονται γι' αυτή τη κάρτα, μιας και είναι πιθανό το πέταγμα του από την τελευταία.
- ACK : Το σήμα αυτό όταν είναι ενεργό σημαίνει ότι το τελευταίο cell που στάλθηκε σε κάποια κάρτα προορισμού έγινε δεκτό από την τελευταία. Σε συνδυασμό με ενεργό το CONGOUT, σημαίνει ότι το cell αν και έγινε δεκτό από την κάρτα προορισμού, προκάλεσε congestion στην τελευταία (σχεδόν γεμάτη ουρά).
- NACK : Το σήμα αυτό όταν είναι ενεργό σημαίνει ότι το τελευταίο cell που στάλθηκε σε κάποια κάρτα προορισμού δεν έγινε δεκτό. Σε συνδυασμό με ενεργό το CONGOUT, σημαίνει ότι δεν έγινε δεκτό λόγω ύπαρξης ισχυρού congestion στην κάρτα προορισμού (γεμάτη ουρά).

Τα παραπάνω σήματα είναι μετάφραση των σημάτων αρνητικής λογικής του CellBus, CBCONG-, ACK- και NACK- αντίστοιχα.

Στο σχήμα 1.3 φαίνεται η περίπτωση μετάδοσης τριών cells σε τρεις κάρτες, μια σχεδόν γεμάτη, μία κανονική και μια πλήρως γεμάτη αντίστοιχα, καθώς και οι τιμές των σημάτων CONGOUT, ACK και NACK που μεταδίδει το CubitPro στον ABRSU. Στην πρώτη περίπτωση το cell γίνεται δεκτό από τη κάρτα προορισμού (ACK = 1) αλλά λόγω της σχεδόν γεμάτης ουράς το CONGOUT = 1. Στη δεύτερη περίπτωση, το cell γίνεται δεκτό (ACK = 1) και αφού η ουρά δεν κινδυνεύει να γεμίσει το CONGOUT = 0. Τέλος, στην τρίτη περίπτωση το cell πετιέται (NACK = 1) λόγω της πλήρως γεμάτης ουράς (CONGOUT = 1).

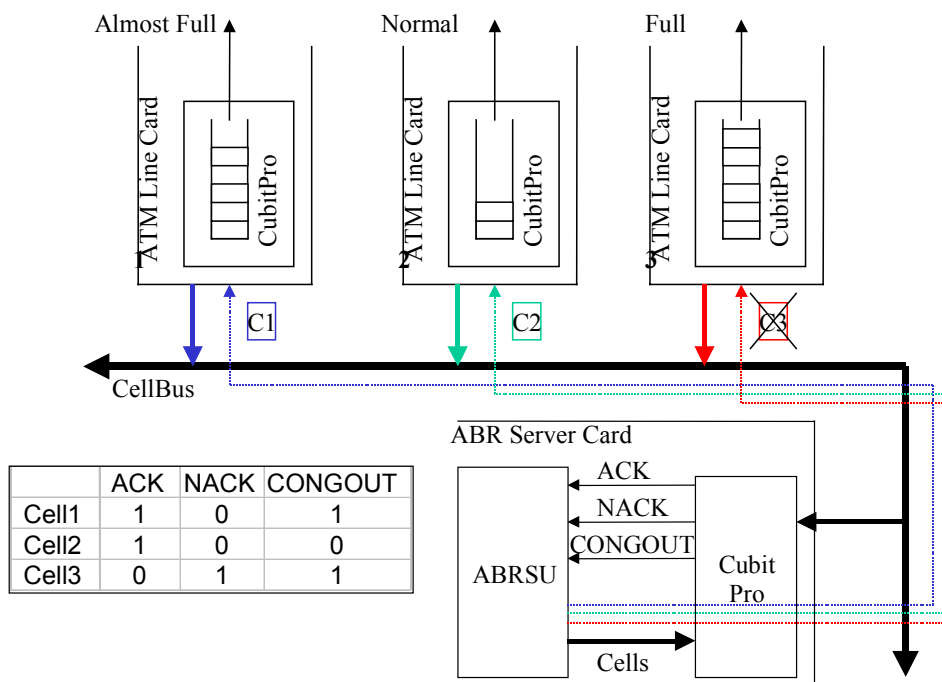


Figure 1-3 :Μηχανισμός Flow Control Κάρτας Εξυπηρετητή ABR

1.5 Η αρχιτεκτονική της Μονάδας Εξυπηρετητή ABR (ABRSU)

Η Μονάδα Εξυπηρετητή ABR υλοποιείται στη Κάρτα Εξυπηρετητή ABR από μία FPGA τύπου EPF 10K200EBC600-1 της Altera. Παρέχει τις παρακάτω λειτουργικότητες στην Κάρτα:

Αποδοχή εισερχόμενης κίνησης τύπου ABR μέσω αφιερωμένης διεπαφής UTOPIA από ολοκληρωμένα CubitPro της κάρτας.

Αναγνώριση των κελιών τύπου RM, ώστε να μπορεί να εφαρμοστεί μαρκάρισμα τύπου RM πάνω τους στην περίπτωση που δεν συμφωνούν με τον τρέχον έλεγχο ροής.

Αποθήκευση των εισερχόμενων κελιών στη Μνήμη SDRAM, με βάση τη ροή στην οποία ανήκουν (per flow queueing). Οι δείκτες των ουρών (Tail Pointers) ενημερώνονται με τις νέες αφίξεις.

Ομαδοποίηση των ροών σε ομάδες ροών (Flow Groups). Η ομαδοποίηση είναι ελεύθερη από κάθε περιορισμό και μπορεί να χρησιμοποιηθεί για την ομαδοποίηση είτε με βάση την κάρτα αποστολής του μετάγωγέα ΔΠΠΟΛΟ, είτε την ποιότητα υπηρεσιών που τους παρέχεται είτε και τις δύο.

Χρονοπρογραμματισμός της αποστολής των αποθηκευμένων κελιών που ανήκουν στην κεφαλή κάθε ουράς ροής. Η μετάδοση βασίζεται στην διαθέσιμη διαμεταγωγή του δίαυλου CellBus. Για αυτή τη λειτουργία, μια μονάδα χρονοπρογραμματισμού (Scheduling block) έχει υλοποιηθεί, που χρησιμοποιεί τις παραμέτρους για την ποιότητα υπηρεσιών για κάθε μία από τις ομάδες ροών για να ζητάει από την ανάγνωση ενός κελιού προς μετάδοση. Αυτό το κελί θα αναγνωσθεί από την μνήμη SDRAM από τον διαχειριστή ουρών που θα ενημερώσει και την κεφαλή της αντίστοιχης ουράς ροής. Όταν η ορθή μετάδοση του κελιού στην κάρτα αποστολής επιβεβαιωθεί τότε ο χώρος στη μνήμη που δέσμευε θα ενταχθεί στη λίστα ελευθέρων και τελικά θα ξαναχρησιμοποιηθεί από άλλο εισερχόμενο κελί.

Μετάδοση των εξερχόμενων κελιών στα ολοκληρωμένα CubitPro της κάρτα μέσω μιας αποκλειστικής διεπαφής εξόδου τύπου UTOPIA. Υπάρχει παράλληλα και η δυνατότητα επαναμετάδοσης ενός κελιού στην περίπτωση που η πρώτη μετάδοση του κελιού πάνω από το CellBus από το CubitPro αποτύχει.

Παροχή διεπαφής μικροεπεξεργαστή για τη διασύνδεση της λογικής της ABRSU με τον μικροεπεξεργαστή MPC860 της κάρτας. Μέσω αυτής της διεπαφής, ο MPC860 μπορεί να αρχικοποιεί ή να τερματίζει ροές – συνδέσει τύπου ABR, που εξυπηρετούνται από την κάρτα. Έχει, παράλληλα, πρόσβαση στις δομές δεδομένων των ουρών καθώς και στις ίδιες τις ουρές που είναι αποθηκευμένες στην μνήμη SDRAM. Ο MPC860 μπορεί επίσης να τροποποιήσει δυναμικά τις παραμέτρους Ελέγχου Ροής για κάθε μία από τις ροές αλλάζοντας τα όρια που χρησιμοποιούνται για το μαρκάρισμα κάθε ροής ξεχωριστά. Μπορεί τέλος να τροποποιήσει και τους χρόνους εξυπηρέτησης που χρησιμοποιεί η Μονάδα Χρονοπρογραμματισμού για να ζητάει την μετάδοση κελιών των ομάδων ροών.

Στο σχήμα 1.4 φαίνεται η εσωτερική οργάνωση της ABRSU.

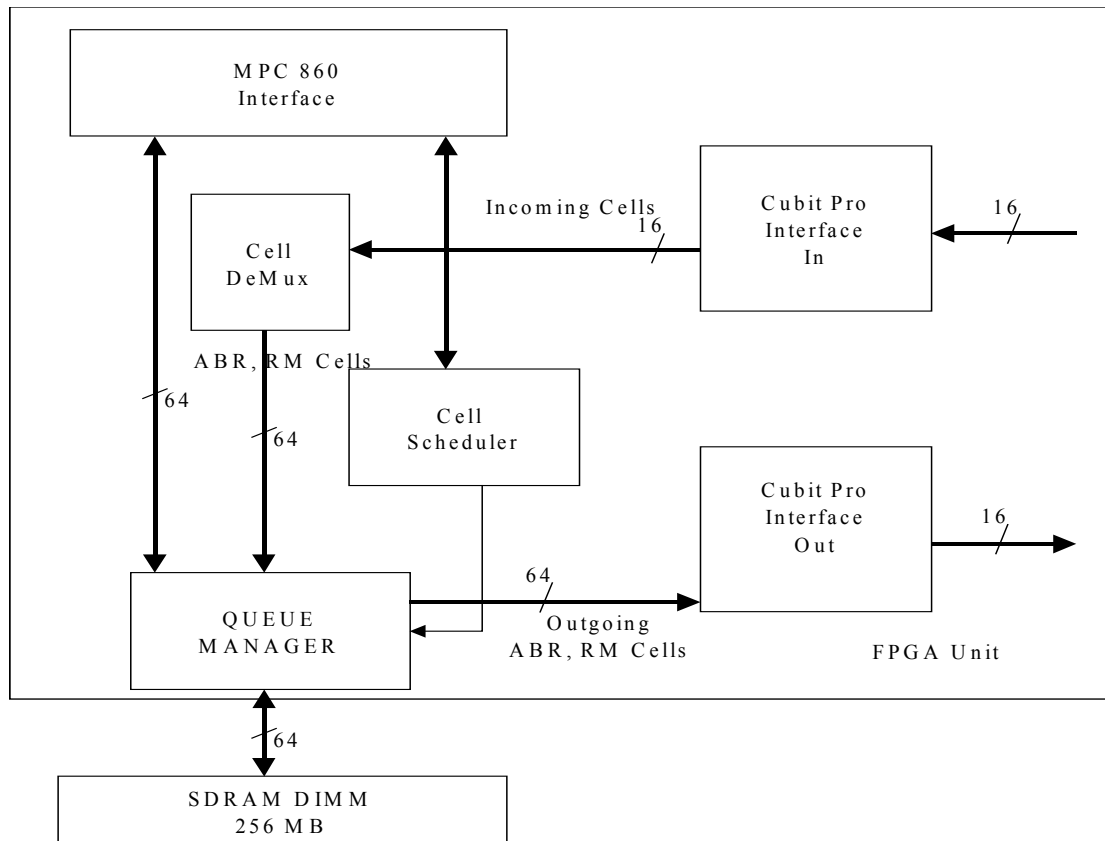


Figure 1-4: Εσωτερικό διάγραμμα του ABRSU

Οι εσωτερικές υπομονάδες της ABRSU παρουσιάζονται στα παρακάτω υποκεφάλαια, ενώ ένα ξεχωριστό κεφάλαιο είναι αφιερωμένο στον διαχειριστή ουρών που είναι και το κύριο θέμα αυτής της εργασίας (Δες Κεφάλαιο 2).

1.5.1 Η Διεπαφή Επεξεργαστή (CPU Interface)

Η υπομονάδα διεπαφής CPU, όπως αναφέρθηκε παραπάνω, είναι υπεύθυνη για την επικοινωνία, αρχικοποίηση και δυναμική τροποποίηση των υπόλοιπων υπομονάδων της ABRSU από τον Επεξεργαστή της Κάρτας. Μέσω αυτής της διεπαφής ο Επεξεργαστής αρχικοποιεί τις δομές δεδομένων των ροών που είναι αποθηκευμένες στην μνήμη SDRAM, θέτοντας συγκεκριμένες εντολές που εκτελούνται από το διαχειριστή ουρών. Αν χρειάζεται, ο διαχειριστής ουρών θα επιστρέψει δεδομένα στον MPC860 μέσω αυτής της διεπαφής. Ο MPC860 μπορεί, επίσης, να τροποποιήσει τις παραμέτρους ποιότητας υπηρεσίας των ομάδων ροών στην υπομονάδα χρονοπρογραμματισμού (Cell Scheduler) [23],[23].

Η εσωτερική οργάνωση του CPU Interface δίνεται στο σχήμα 1.5. Μιας και ο διάυλος του MPC860 στην κάρτα δουλεύει σε συχνότητες ρολογιού που κυμαίνονται από 10 μέχρι 25 MHz ενώ η FPGA μπορεί να πετύχει υψηλότερες συχνότητες εσωτερικού ρολογιού, η Διεπαφή CPU πρέπει να χρησιμοποιεί και τις δύο διαφορετικές συχνότητες ρολογιού καθώς και κάποιο κύκλωμα συγχρονισμού.

Από το μέρος του διαύλου του επεξεργαστή η μονάδα διεπαφής χρησιμοποιεί τα pins του διαύλου σαν είσοδο/έξοδο. Αυτά είναι 32 για τα δεδομένα, 22 για διευθυνσιοδότηση και μερικά επιπλέον σήματα ελέγχου signals (chip select,

Read/Write etc). Μέσω αυτών των pins ο Μικροελεγκτής μπορεί να γράφει και σε συγκεκριμένους καταχωρητές και έτσι να διαμορφώνει δυναμικά την λειτουργία και τις δομές των υπομονάδων Queue Manager Block και Cell Scheduler Block. Αντίστοιχα μπορεί να λάβει δεδομένα εκτελώντας μια πράξη ανάγνωσης σε καταχωρητή που ενημερώνεται από την υπομονάδα Cell Scheduler Block.

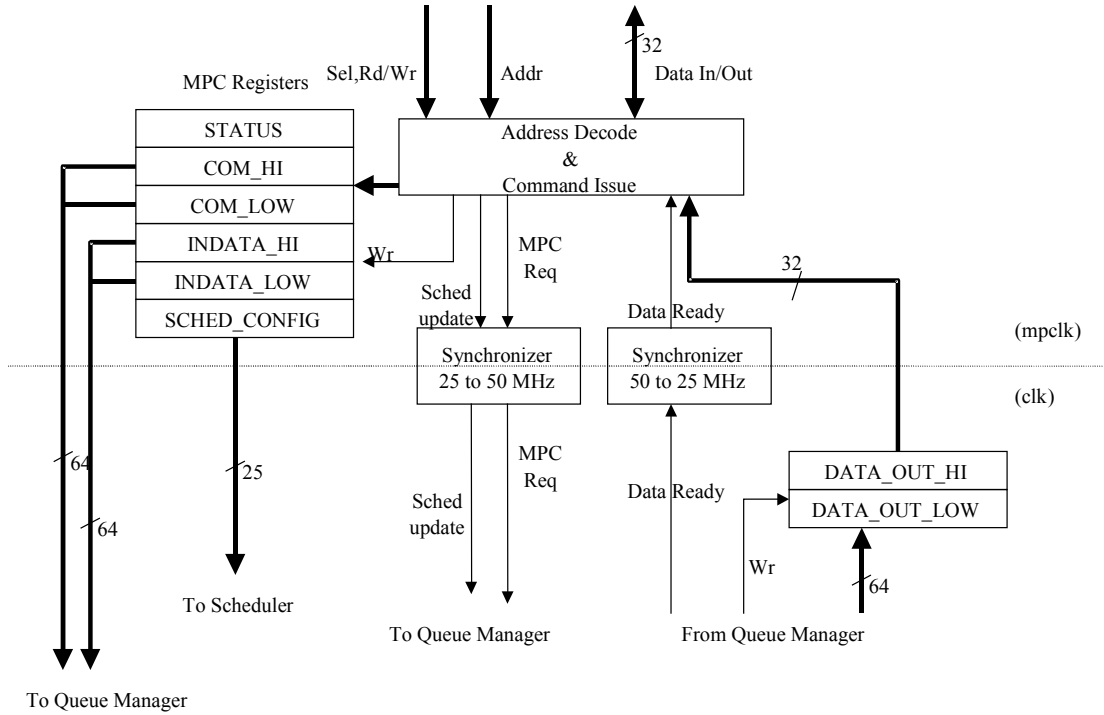


Figure 1-5: Διάγραμμα της διεπαφής CPU

Υπεύθυνη για την εγγραφή των καταχωρητών με βάση τα pins του διαύλου είναι ένα κομμάτι λογικής που στο σχήμα φαίνεται με το όνομα Address Decode and Command Issue. Η λειτουργία της είναι να αποκωδικοποιεί το περιεχόμενο του διαύλου διευθύνσεων και να εκτελεί ανάγνωση ή εγγραφή σε αυτούς με βάση το Rd/Wr. Ακόμη η εγγραφή σε ένα συγκεκριμένο καταχωρητή από τον Μικροελεγκτή έχει σαν αποτέλεσμα την ενεργοποίηση του σήματος Mpc_Req που αφού περάσει από κύκλωμα συγχρονισμού από τα 25 MHz στα 50 MHz θα ειδοποιήσει την υπομονάδα Queue Manager ότι υπάρχει μια εντολή του Μικροελεγκτή προς εκτέλεση. Αντίστοιχη λειτουργία έχει και το σήμα Sched Update που αφού περάσει από κύκλωμα συγχρονισμού ειδοποιεί την υπομονάδα Cell Scheduler ότι ο αντίστοιχος καταχωρητής που παρέχει δεδομένα διαμόρφωσης της τελευταίας έχει γεμίσει με δεδομένα.

1.5.2 Μονάδα διεπαφής CPU – διεπαφή με Διαχειριστή Ουρών

Η υπομονάδα MPC 860 Interface Block μπορεί να διαμορφώσει την Μονάδα αποθήκευσης Cells και Δομών Ουρών (SDRAM) που ελέγχεται από την υπομονάδα Queue Manager καθώς και με τις δομές που διατηρεί η τελευταία στην ίδια μνήμη με την θέση συγκεκριμένων εκτολών που του παρέχονται. Κάποιες από τις εντολές συνοδεύονται από δεδομένα ενώ άλλες ενώ άλλες όχι ενώ κάποιες επιστρέφουν δεδομένα.

Οι εντολές που μπορεί να θέσει η CPU στην Διεπαφή CPU που αφορούν τον Διχειριστή Ουρών και οι καταχωρητές που χρησιμοποιούνται από κάθε μία είναι οι εξείς:

- 1) Open Flow (Ανοιξε Ροή) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW.
- 2) Close Flow (Κλείσε Ροή) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW.
- 3) Write (Γράψε) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW, INDATA_HI, INDATA_LOW.
- 4) Read (Διάβασε) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW.
- 5) Read Counter (Διάβασε Μετρητή) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW.
- 6) Change Parameters (Άλλαξε Παραμέτρους) : Απαιτεί εγγραφή καταχωρητών COM_HI, COM_LOW.

Για κάθε μια από τις πιθανές εντολές, πρέπει να γραφτούν 2 τουλάχιστον καταχωρητές. Είναι απαραίτητο ο τελευταίος καταχωρητής που θα εγγραφεί να είναι ο COM_LOW, μιας και η εγγραφή του θέτει το σήμα Mpc_Req που αφού περάσει από κύκλωμα συγχρονισμού από τα 25 MHz στα 50 MHz θα ειδοποιήσει την υπομονάδα Queue Manager ότι υπάρχει μια εντολή του Μικροελεγκτή προς εκτέλεση.

Οι εντολές Read και Read Counter είναι εντολές ανάγνωσης και αναμένουν επιστροφή δεδομένων από την υπομονάδα Queue Manager. Όταν αυτή ετοιμάσει τα δεδομένα τα αποστέλλει στους καταχωρητές DATA_OUT_HI και DATA_OUT_LOW. Η εγγραφή των καταχωρητών αυτών από την υπομονάδα Queue Manager έχει σαν αποτέλεσμα την θέση του σήματος Data-Ready που αφού περάσει από κύκλωμα συγχρονισμού από τα 50 MHz στα 25 MHz θα θέσει αντίστοιχο bit στον καταχωρητή STATUS. Ο Μικροελεγκτής εξετάζοντας την τιμή αυτού του bit, μαθαίνει ότι τα δεδομένα είναι έτοιμα για ανάγνωση.

1.5.3 Μονάδα διεπαφής CPU - διεπαφή με Cell Scheduler

Η υπομονάδα MPC 860 Interface δίνει την δυνατότητα στον Μικροελεγκτή να διαμορφώσει την ποιότητα της εξυπηρέτησης των 64 Ομάδων Ροών που υποστηρίζει η Μονάδα Εξυπηρέτησης ABR μέσω της υπομονάδας Cell Scheduler. Η διαμόρφωσή γίνεται με την εγγραφή του καταχωρητή SCHED_CONFIG. Σε αυτόν ο Μικροελεγκτής γράφει τιμές που θα ενημερώσουν την μνήμη που διατηρεί για την ποιότητα υπηρεσίας της κάθε ομάδας ροής η υπομονάδα Cell Scheduler. Η εγγραφή αυτού του καταχωρητή έχει σαν αποτέλεσμα τη θέση του σήματος Sched Update που αφού περάσει από κύκλωμα συγχρονισμού από τα 25 MHz στα 50 MHz θα ειδοποιήσει την υπομονάδα Cell Scheduler ότι υπάρχει νέα τιμή για κάποια Ομάδα Ροών προς ενημέρωση του τελευταίου.

1.5.4 Η υπομονάδα Cell Demultiplexor

Η υπομονάδα Cell Demultiplexor είναι υπεύθυνη για τη συγκέντρωση ενός ολόκληρου κελιού σε 7 λέξεις των 64 bit η κάθε μία, ώστε να αποσταλούν στο διαχειριστή ουρών για αποθήκευση στην κατάλληλη ουρά. Αυτή η υπομονάδα είναι απαραίτητη μιας και τα κελιά πρώτα μπαίνουν στην μονάδα διεπαφής CubitPro και αποθηκεύονται σε μία FIFO των 16 bit και πρέπει να οργανωθούν σε μια fifo των 64 bit ώστε η αποθήκευσή τους στην SDRAM να γίνει σε αλληπαλούς κύκλους των 64 bity όπως απαιτεί η διεπαφή με την SDRAM.

Η υπομονάδα εξετάζει το σήμα avail που έρχεται από τη μονάδα διεπαφής CubitPro για να δει αν υπάρχει διαθέσιμο κελί στην εσωτερική fifo της τελευταίας. Αν το σήμα είναι θετικό προχωράει στο γέμισμα ενός καταχωρητή μεγέθους 64 bit με τις λέξεις των 16 bit. Όταν ο καταχωρητής γεμίσει μια εγγραφή γίνεται από την FSM στην fifo των 64 bit με τα δεδομένα του καταχωρητή. Αυτή η fifo είναι 7x64. Όταν γεμίσει, ένα κελί είναι πλέον έτοιμο να αποθηκευτεί στην SDRAM και το σήμα fifo_full χρησιμοποιείται σαν request στο διαχειριστή ουρών για αποθήκευσή του. Κανένα νέο κελί δεν μπορεί να εισαχθεί από την FSM αν η fifo δεν αδειάσει πρώτα. Αυτό γίνεται με αλληπαλά διαβάσματα από το διαχειριστή ουρών από τη fifo. Η πρώτη λέξη κάθε κελιού περιέχει και το ID (ταυτότητα) της ροής στην οποία ανήκει αυτό.

Η οργάνωση της υπομονάδας Cell Demultiplexor φαίνεται στο σχήμα 1.6

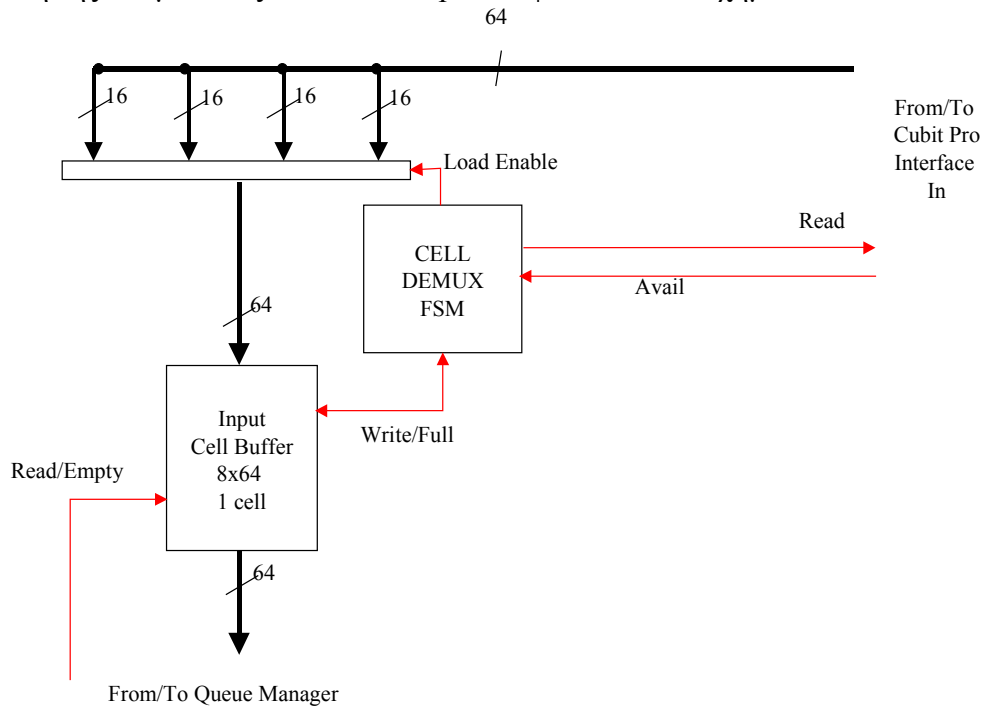


Figure 1-6: Εσωτερικό διάγραμμα του Cell Demultiplexor

1.5.5 Διεπαφές UTOPIA

Υπάρχουν δύο υπομονάδες διεπαφής UTOPIA, και οι δύο των 16 bit, μέσα στην ABRSU. Η πρώτη ονομάζεται UTOPIA Input Interface και χρησιμοποιείται για να εισάγει κελιά από τα ολοκληρωμένα CubitPro στην ABRSU. Η δεύτερη ονομάζεται UTOPIA Output Interface και χρησιμοποιείται για την αποστολή των εξεχόμενων

από την ABRSU κελιών στα ολοκληρωμένα Cubit Pro προς μετάδοση πάνω από το CellBus προς την κάρτα προορισμού τους.

1.5.6 Ο Προγραμματιστής Cells (Cell Scheduler) [23]

Η Μονάδα Cell Scheduler του ABRSU είναι υπεύθυνη για τον χρονοπρογραμματισμό της μετάδοσης, πάνω από το CellBus, των cells που βρίσκονται αποθηκευμένα στις ουρές των εικονικών συνδέσεων που εξυπηρετεί η Κάρτα Εξυπηρετητή ABR. Αυτός ο χρονοπρογραμματισμός γίνεται με τέτοιο τρόπο ώστε να μοιράζεται δίκαια το διαθέσιμο bandwidth, για την Κάρτα Εξυπηρετητή ABR, bandwidth.

Ο Cell Scheduler μαθαίνει μέσω πληροφορίας Flow Control που δέχεται από τις Συσκευές Μεταγωγής CubitPro, ποιες από τις κάρτες προορισμού είναι διαθέσιμες για λήψη Cell (δηλαδή δεν υποφέρουν από συμφόρηση) και έτσι περιορίζει το πλήθος των ουρών που μπορούν να μεταδώσουν ένα cell σε αυτές που α) έχουν cell για να στείλουν και β) η κάρτα προορισμού έχει χώρο για να δεχθεί το cell. Στην συνέχεια , αφού επιλεγούν οι ικανές (σύμφωνα με τα παραπάνω) ουρές, αυτές εξυπηρετούνται με πολιτική Round Robin, ώστε να μοιραστεί δίκαια το διαθέσιμο bandwidth και να υπάρχει ένα άνω όριο στην καθυστέρηση κάθε cell στην ουρά μιας εικονικής σύνδεσης που δεν μεταδίδει παραπάνω από το επιτρεπτό throughput.

1.5.7 Ο Διαχειριστής Ουρών (Queue Manager)

Ο Διαχειριστής Ουρών, τμήμα του ABRSU, δέχεται cells, στην εισερχόμενη κατεύθυνση, από τα CubitPro. Τα cells περιέχουν το πεδίο των 16 bit που καθορίζει την ουρά που θα αποθηκευτούν καθώς και τους headers που καθορίζουν σε ποια κάρτα προορισμού θα σταλθεί το cell όταν έρθει η ώρα να μεταδοθεί. Τα cell αποθηκεύονται στις αντίστοιχες ουρές με πολιτική FIFO. Μια μονάδα μνήμης SDRAM χρησιμοποιείται εδώ για αποθήκευση πληροφορίες διαχείρισης των ουρών (head και tail pointers), καθώς και για την αποθήκευση των ίδιων των Cell.

Ο Queue Manager απαντά στις αιτήσεις του Cell Scheduler βγάζει από τις ουρές cells για μετάδοση. Ακόμη απαντά στις αιτήσεις του επεξεργαστή MPC860 για την δημιουργία νέων ουρών VP/VC (CAC – Call Admission Control για μια νέα σύνδεση ABR) με την δέσμευση ενός από το πεπερασμένο πλήθος ελεύθερων (μη χρησιμοποιούμενων από τις 64 χιλιάδες) ουρών.

Τέλος ο διαχειριστής ουρών διαχειρίζεται και τον ελεύθερο χώρο στην μνήμη αποθήκευσης δεδομένων μέσω μίας ουράς ελεύθερου χώρου (free-list) και ενημερώνει την δομή αυτή κάθε φορά που ένα cell αποθηκεύεται ή μεταδίδεται.

2 Η μονάδα Διαχειριστή Ουρών (Queue Manager)

Σε αυτό το κεφάλαιο παρουσιάζουμε τον Διαχειριστή Ουρών. Αυτή η λειτουργική μονάδα υλοποιήθηκε με τη χρήση της γλώσσας Verilog και συνθέθηκε με το εργαλείο σύνθεσης MaxPlus II, για να χωρέσει σε μια FPGA της Altera. Χρησιμοποιεί μία μόνο μονάδα μνήμης SDRAM DIMM για την αποθήκευση των εισερχόμενων κελιών τύπου ABR, καθώς και για τα δεδομένα για τις λογικές ουρές των κελιών. Έτσι μειώνεται το πλήθος των χρησιμοποιούμενων pins και των συρμάτων πάνω στη κάρτα και κατά συνέπεια το συνολικό κόστος της κάρτας. Η αρχιτεκτονική αποθήκευσης που επιλέχθηκε είναι αυτή της ανά-ροής αποθήκευσης (per-flow queueing). Κάθε ξεχωριστή ροή δηλαδή που εξυπηρετείται από την κάρτα δεσμεύει μία ξεχωριστή ουρά για την αποθήκευση των κελιών της. Το μέγιστο πλήθος των υποστηριζόμενων ουρών κάθε στιγμή είναι 64 χιλιάδες, αριθμός υπεραρκετός για ένα μεταγωγέα άκρης. Μια λίστα από ελεύθερους αποθηκευτικούς χώρους διατηρείται επίσης, για να διατηρεί όλους τους χώρους κελιών στη μνήμη που δεν χρησιμοποιούνται. Κατά συνέπεια υλοποιείται δυναμική παραχώρηση μνήμης (dynamic memory allocation), επιτρέποντας στις ροές να έχουν μεταβλητό μέγιστο μέγεθος για καλύτερη χρησιμοποίηση της μνήμης. Οι ροές επίσης οργανώνονται σε ομάδες ροών (Flow Groups) με την χρήση κυκλικών λιστών. Δεδομένα αυτών των κυκλικών λιστών βρίσκονται σε μια εσωτερική στην FPGA μνήμη τύπου SRAM. Η αρχικοποίηση των ροών και ο καθορισμός της ομάδας ροών στην οποία ανήκουν, γίνεται από τον μικροεπεξεργαστή της κάρτας κατά την έναρξη των συνδέσεων, με ειδικές εντολές αρχικοποίησης που θέτει ο τελευταίος μέσω της μονάδας διεπαφής CPU. Οι παράμετροι του ελέγχου ροής μπορούν επίσης να τεθούν από τον επεξεργαστή, με τον καθορισμό του μέγιστου μεγέθους για κάθε ουρά ξεχωριστά. Αν μια ουρά ξεπεράσει το όριό της το μαρκάρισα τυπου RM και EFCI, ξεκινάει, όπως αυτό περιγράφεται στο ATM Forum.

Μιας και μόνο μία μνήμη χρησιμοποιείται για τις ανάγκες του Διαχειριστή Ουρών, δεν είναι δυνατόν να γίνουν παράλληλες προσπελάσεις σε μνήμες. Έτσι μόνο σειριακές προσπελάσεις είναι δυνατές κατά την αποθήκευση και την έξοδο ενός κελιού από μία ουρά. Αυτό το γεγονός κάνει τη μνήμη τροχοπέδη στην ικανότητα διαμεταγωγής του συστήματος. Προκειμένου να μειωθούν οι προσπελάσεις στη μνήμη, υλοποιήθηκε προ-ανάθεση χώρου (buffer preallocation) σε κάθε μία από τις 64K ουρές καθώς και παράκαμψη λίστας ελευθέρων (free list bypassing) με τη βοήθεια εξωτερικής του Διαχειριστή λογικής (του χρονοπρογραμματιστή κελιών). Ο ελεγκτής της SDRAM που είναι μέρος του διαχειριστή ουρών για να χειρίζεται τις προσπελάσεις στην SDRAM είναι επίσης προσεκτικά σχεδιασμένος ώστε να κάνει αλληπάλληλες προσβάσεις χωρίς την απώλεια κύκλων ρολογιού. Οι εντολές enqueue (αποθήκευση κελιού σε μια ουρά) και dequeue (έξοδος ενός κελιού από μία ουρά) είναι επίσης σχεδιασμένες ώστε να μην χάνονται κύκλοι λόγω αλληπάλληλων προσβάσεων στην μνήμη με εξαρτημένα δεδομένα αναμεταξύ τους.

Εκτός από την διεπαφή με τον επεξεργαστή και την SDRAM, ο διαχειριστής ουρών διατηρεί μια διεπαφή για την είσοδο των κελιών προς αποθήκευση και μια διεπαφή με το υλικό χρονοπρογραμματισμού που ζητάει την έξοδο των κελιών συγκεκριμένων ομάδων ροών προς μετάδοση. Η αίτηση εξυπηρετείται από τον Διαχειριστή που διαβάσει ένα κελί μιας συγκεκριμένης ομάδας ροών και το στέλνει στη διεπαφή σε λέξεις των 64 bit.

2.1 Η αρχιτεκτονική του Διαχειριστή Ουρών

Στο σχήμα 2.1 φαίνεται η εσωτερική οργάνωση του Διαχειριστή Ουρών.

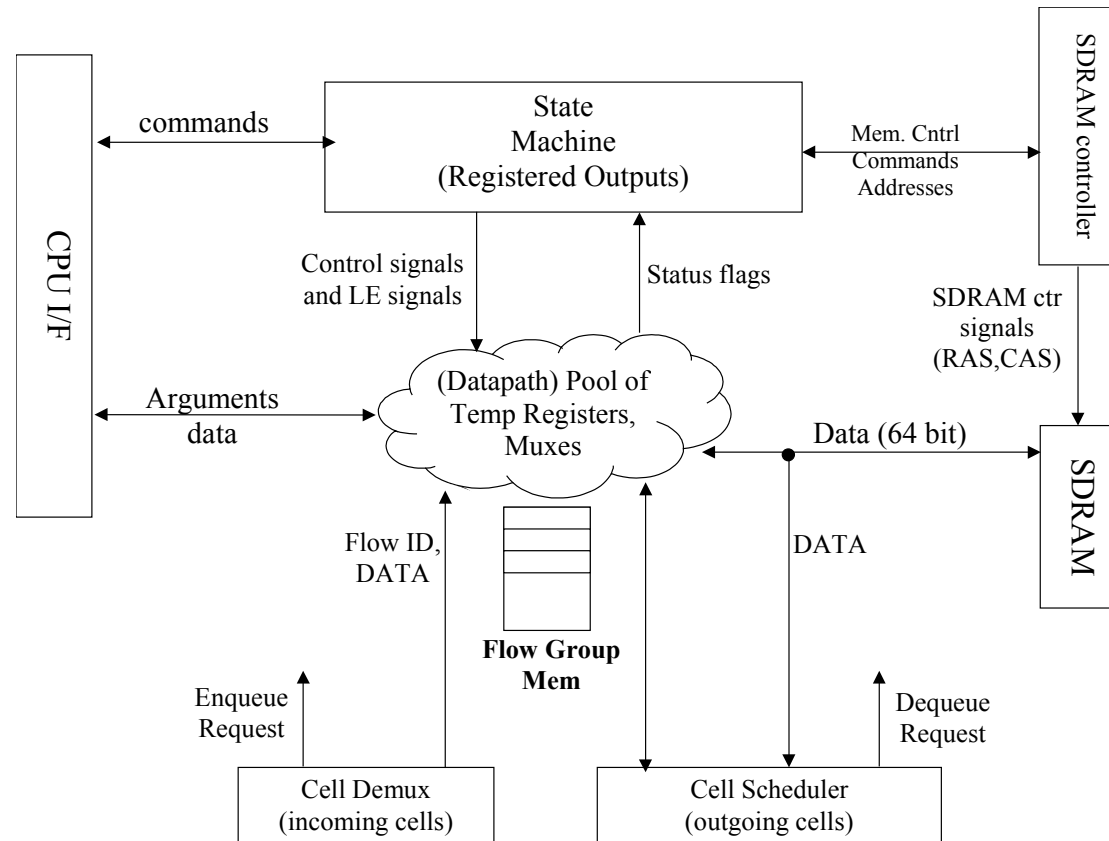


Figure 2-1: Το εσωτερικό διάγραμμα του διαχειριστή ουρών

Όπως φαίνεται στο σχήμα ο Διαχειριστής Ουρών αποτελείται από τα παρακάτω:

- **Μηχανή Καταστάσεων (State Machine):** Αυτή η μηχανή πεπερασμένων καταστάσεων ουσιαστικά αποτελείται από πολλές μηχανές, κάθε μία από τις οποίες είναι αφιερωμένη στον έλεγχο της υπόλοιπης λογικής καταχωρητών και μνήμης του διαχειριστή, με σκοπό την υλοποίηση μιας από τις εντολές του διαχειριστή. Όλες αυτές οι μηχανές καταστάσεων ενεργοποιούνται από την μηχανή καταστάσεων TOP που αποφασίζει ποιά εντολή πρέπει να εκτελεστεί μία δεδομένη στιγμή.
- **Μονοπάτι δεδομένων (Datapath):** Πρόκειται για προσωρινούς καταχωρητές που φορτώνονται με δεδομένα που προέρχονται από όλες τις λειτουργικές μονάδες που επικοινωνούν με τον Διαχειριστή ουρών. Φορτώνονται με βάση σήματα ελέγχου που προέρχονται από τη μηχανή πεπερασμένων καταστάσεων και έπονται κάποιων πολυπλεκτών που επιλεγουν την προέλευση των δεδομένων που αποθηκεύονται. Η μηχανή καταστάσεων και πάλι ελέγχει με σήματα αυτούς τους πολυπλέκτες. Ονόματα τέτοιων καταχωρητήν και των αντίστοιχων πολυπλεκτών τους είναι : Flow ID, Head Pointer, Tail Pointer, Cell Counter, Hi Watermark, Flow Group ID (δείκτες μιας δεδομένης ουράς που ενημερώνεται από την

εκάστοτε εντολή), Free List Head, Free List Tail, Free List Counter (Πληροφορία για τη λίστα ελευθέρων χώρων της μνήμης κελιών SDRAM) κλπ.

- **Μνήμη Ομάδων Ροών (Flow Group Memory)** : Αυτή η μνήμη κρατά δεδομένα για τις 64 ομάδες ροών που υποστηρίζονται από το Διαχειριστή Ροών. Είναι μνήμη 64 λέξεων και αποθηκεύει την κεφαλή της κυκλικής λίστας των ενεργών ροών που ανήκουν σε κάθε ροή, την ουρά της λίστας, και την κατάσταση της κάθε ομάδας (αν η Ομάδα ροών έχει κάποια ενεργή ροή ή όχι). Μια ροή θεωρείται ενεργή όταν υπάρχει ένα κελί αυτής αποθηκευμένο στη μνήμη κελιών SDRAM που θα πρέπει να μεταδοθεί κάποια στιγμή στο μέλλον. Η μηχανή καταστάσεων ελέγχει αυτή τη μνήμη.
- **Ο ελεγχτής SDRAM (The SDRAM controller)**: Αυτή η υπομονάδα υλοποιεί τις προσπελάσεις στην μνήμη SDRAM DIMM εκ μέρους της μηχανής καταστάσεων. Θέτει τα pin ελέγχου της SDRAM DIMM και την προγραμματίζει να γράψει ή να διαβάσει 1,2,8 λέξεις των 64 bit. Τα δεδομένα δίνονται από την υπόλοιπη λογική του διαχειριστή κατευθείαν, σε συγχρονισμό με τα σήματα του ελεγκτή. Δεν απαιτείται συγχρονισμός μιας και η FPGA που υλοποιεί τον διαχειριστή ουρών και η SDRAM DIMM χρησιμοποιούν το ίδιο, προσεκτικά καταναμημένο ρολόι.

Ο Διαχειριστής ουρών διατηρεί διεπαφή με τις παρακάτω υπόμονάδες:

- **Την διεπαφή CPU (The CPU Interface)**: Αυτή η διεπαφή επιτρέπει μια CPU να μπορεί να διαμορφώσει τον Διαχειριστή Ουρών , να αρχικοποιήσει ροές, να διαμορφώσει τις παραμέτρους ροής και να αντλήσει δεδομένα που βοηθούν στην αφαίρεση λαθών στο σύστημα.
- **Η διεπαφή με τον Cell Demultiplexor) The Cell Demultiplexor Interface**: Μέσω αυτής της διεπαφής ο διαχειριστής μπορεί να προσπελάσει την fifo 8x64 που περιέχει το επόμενο κελί που πρέπει να εισαχθεί στη μνήμη. Αν η fifo είναι γεμάτη, το σήμα `fifo_full` χρησιμοποιείται ως αίτηση `enqueue` στην μηχανή καταστάσεων.
- **Η διεπαφή Χρονοπρογραμματιστή κελιών (The Cell Scheduler interface)** : Αυτή είναι η διεπαφή με την μονάδα χρονοπρογραμματισμού που υλοποιεί την χρονοπρογραμματιστική αρχιτεκτονική του συστήματος. Μία αίτηση `dequeue` μπορεί να δοθεί στο διαχειριστή μέσα από αυτή την διεπαφή, μαζί με τον κωδικό της ομάδας ροής από την οποία πρέπει να προέρχεται το κελί που πρέπει να εξέλθει. Ο διαχειριστής απαντά στέλνοντας το κελί σε λέξεις των 64 bit μαζί με την διεύθυνση στην μνήμη SDRAM που βρίσκεται το κελί που εξέρχεται. Όταν υλοποιείται η παράκαμψη λίστας ελευθέρων (Free list bypassing), αυτές οι διευθύνσεις δίνονται πίσω στο διαχειριστή ουρών για την αποθήκευση νεοαφιχθέντων κελιών.
- **Η διεπαφή με την μνήμη SDRAM DIMM (The SDRAM DIMM interface)**: Αυτή η διεπαφή ελέγχεται όπως περιγράφεται παραπάνω από τον ελεγκτή SDRAM και μόνο ο δίαυλος δεδομένων της μνήμης ελέγχεται από την υπόλοιπη λογική του διαχειριστή ουρών.

2.2 Ροές

Η Ροή είναι η βασική δομή δεδομένων που υποστηρίζεται από το Διαχειριστή Ουρών. Η σημασία της είναι αυτή της ροής κίνησης μιας συγκεκριμένης δικτυακής

σύνδεσης που εξυπηρετείται από τον διαχειριστή ουρών. Ο διαχειριστής υποστηρίζει 64K ταυτόχρονες συνδέσεις. Για να τις διαχωρίσει αναμεταξύ τους, σε κάθε μία από τις ροές ανατίθεται ένας συγκεκριμένος αριθμός ταυτότητας (ID number) που δεσμεύεται από την CPU κατά την αρχικοποίηση της σύνδεσης – ροής. Μιας και το μέγιστο πλήθος των υποστηριζόμενων ροών είναι 64K, το ID πρέπει να είναι ένας αριθμός πλάτους 16 bit.

Πληροφορία για κάθε ροή βρίσκεται αποθηκευμένη στις Εγγραφές Ροών (Flow Records). Αυτές οι εγγραφές δεσμεύουν πάντα χώρο μέσα στην SDRAM. Η δομή δεδομένων που αντιπροσωπεύει μία ροή είναι αυτή της σειριακής λίστας μονής διασύνδεσης που φαίνεται στο σχήμα 2.2. Η εγγραφή κάθε ροής περιέχει πληροφορία για αυτή τη λίστα όπως, τον δείκτη κεφαλής, το δείκτη τέλους, τον μετρητή μεγέθους (σε κελιά) της λίστας, κάποια bits κατάσταση, και τις παραμέτρους ελέγχου ροής. Μια λίστα θεωρείται χρησιμοποιούμενη όταν ο επεξεργαστής της έχει αναθέσει μια εισερχόμενη ροή και ενεργή όταν υπάρχει έστω και ένα κελί αυτής αποθήκευμένο μέσα στην SDRAM. Η εγγραφή ροής παρουσιάζεται εκτενώς στο υποκεφάλαιο 2.5.

Όταν ένα κελί που ανήκει σε μία ροή (που καθορίζει το ID που κουβαλάει η επικεφαλίδα του) φτάσει στο Διαχειριστή ουρών, αυτό γίνεται enqueue στη λίστα. Ο διαχειριστής το κάνει αυτό, με το να γράφει το κελί στο ελεύθερο χώρο στο τέλος της αντίστοιχης λίστας και να γράφει τον δείκτη next_pointer αυτού του χώρου να δείχνει σε ένα νέο ελεύθερο χώρο (πρέπει πάντα να υπάρχει ένας ελεύθερος χώρος στην ουρά της λίστας, για λόγους που εξηγούνται στο υποκεφάλαιο 2.11). Έλευθεροι αποθηκευτικοί χώροι για αυτό το σκοπό παίρνονται από τη λίστα ελευθέρων (Free List) που διατηρείται από το Διαχειριστή Ουρών. Ο δείκτης ουράς (Tail Pointer) της Εγγραφής Ροής τείθεται να δείχνει στο καινούργιο ελεύθερο αποθηκευτικό χώρο, ενώ και ο μετρητής αυξάνεται επίσης.

Όταν ένα κελί μιας ροής πρέπει να γίνει dequeue και να σταλθεί στο χρονοπρογραμματιστή κελιών για αποστολή πάνω από το cellbus, ο διαχειριστής ουρών διαβάζει το κελί που είναι η κεφαλή της αντίστοιχης λίστας από τη μνήμη και θέτει το δείκτη κορυφής της εγγραφής ροής να δείχνει στον επόμενο χώρο. Αν η μετάδοση του κελιού επιτύχει τότε ο αποθηκευτικός του χώρος μπορεί πλέον να εισαχθεί στη λίστα ελευθέρων ή να ξαναχρησιμοποιηθεί ως ουρά σε μια πράξη enqueue (Δες παράκαμψη λίστας ελευθέρων υποκεφάλαιο 2.10). Το σχήμα 2.2 δίνει την δομή μιας ροής, μια πράξη enqueue και μια dequeue.

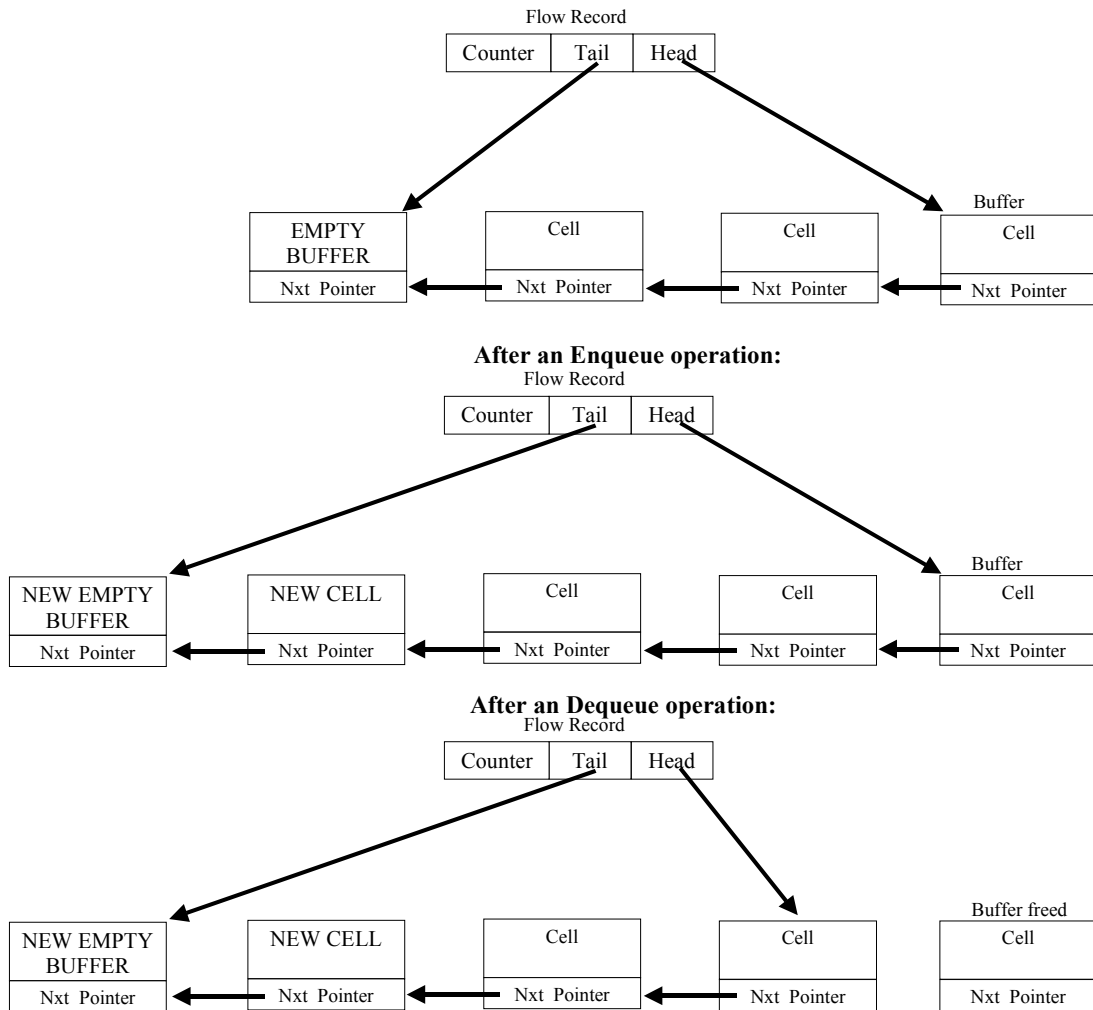


Figure 2-2: Η δομή των ροών και η αλλαγές σε αυτή μετά από μία εντολή enqueue και μία εντολή dequeue.

2.3 Ομάδες Ροών (Flow Groups)

Οι 64 χιλιάδες ροές που υποστηρίζονται από το Διαχειριστή Ουρών οργανώνονται σε υψηλότερό επίπεδο σε ομάδες ροών. Όλες οι ροές που χρησιμοποιούνται από συνδέσεις πρέπει να ανήκουν σε μία ομάδα ροής. Το πλήθος των ομάδων ροών που υποστηρίζονται από το διαχειριστή ουρών είναι 64. Ο λόγος αυτής της ομαδοποίησης είναι η υποστήριξη της χρονοπρογραμματιστικής ομάδας στο χρονοπρογραμματισμό της εξερχόμενης κίνησης τύπου ABR. Ο ανά ροή χρονοπρογραμματισμός με βάση παραμέτρους ποιότητας εξυπηρέτησης (QoS) είναι πολύ δύσκολος στην υλοποίηση του για μεγάλο αριθμό ροών. Για αυτό το λόγο οι ροές οργανώνονται σε ομάδες από το διαχειριστή ουρών, και κάθε ομάδα ροών έχει τις δικές της παραμέτρους QoS που διατηρούνται στο Χρονοπρογραμματιστή που παρέχει στις ομάδες διαμεταγωγή με βάση αυτές τις παραμέτρους. Οι ροές κάθε ομάδας παίρνουν ίδια διαμεταγωγή μιάς και ο Διαχειριστής ουρών τις εξυπηρετεί κυκλικά (Round Robin).

Η δομή δεδομένων που υλοποιεί μία Ομάδα Ροών είναι αυτή της διπλά συνδεδεμένης κυκλικής λίστας που φαίνεται στο σχήμα 2.3. Η κεφαλή και η ουρά της κάθε λίστας (υπάρχουν 64 τέτοιες λίστες) είναι αποθηκευμένες στην μνήμη Ομάδων Ροών (Flow group memory). Μόνο ενεργές ροές υπάρχουν σε αυτές τις λίστες. Αν μια ροή μετα από μια πράξη dequeue γίνεται ανενεργή τότε αφαιρείται από την κυκλική λίστα. Όταν ένα νέο κελί γίνει φτάσει στη ροή και ξαναγίνει ενεργή, επαναεισάγεται στη λίστα ως ουρά.

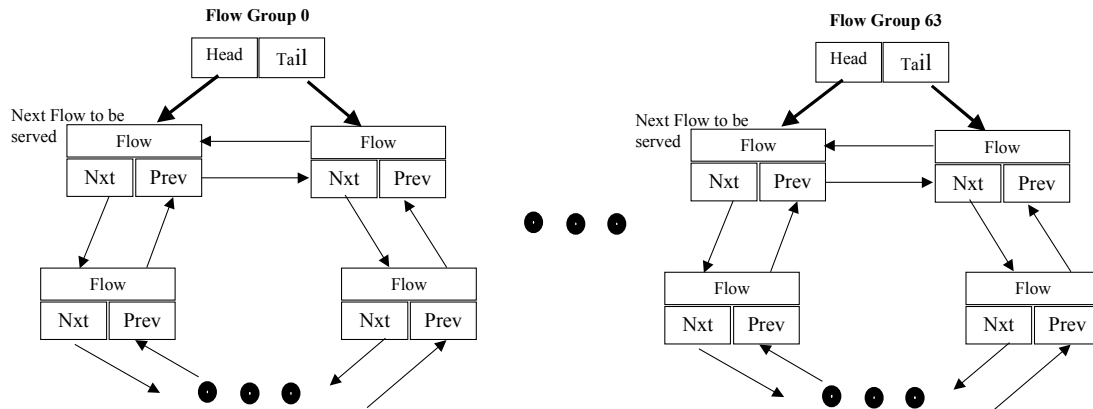


Figure 2-3: Οι 64 κυκλικές λίστες των ομάδων ροών

Αν ο Χρονοπρογραμματιστής κελιών ζητήσει την πράξη dequeue από μια ομάδα ροών, η ροή που είναι κεφαλή στη κυκλική λίστα θα επιλεγεί για την αποστολή του πρώτου της κελιού και η ροή θα γίνει απο κεφαλή ουρά της λίστας ενώ η επόμενη στη σειρά ροή θα γίνει κεφαλή.

Η κυκλική λίστα είναι διπλάσυνδεδεμένη μιας και μια σύνδεση που αντιπροσωπεύεται από μία ροή μπορεί να κλείσει. Σε αυτή τη περίπτωση η ροή πρέπει να αφαιρεθεί από τη λίστα και η τελευταία να παραμείνει συνδεδεμένη με πράξεις χρόνου $O(1)$. Οι δείκτες προς την επόμενη και τη προηγούμενη ροή κρατούνται στην εγγραφή της κάθε ροής.

2.4 Οι εντολές του Διαχειριστή ουρών

Ο Διαχειριστής ουρών μπορεί να δεχθεί ένα εύρος εντολών και να τις εκτελέσει ενεργοποιώντας μία μηχανή καταστάσεων για κάθε μία από αυτές, για να ελέγξει τις δομές δεδομένων, τους καταχωρητές και τις μνήμες. Οι πιο σημαντικές εντολές είναι οι Enqueue και Dequeue εντολές που θέτονται από τον Cell Demux και το Cell Scheduler (Χρονοπρογραμματιστής Κελιών) αντίστοιχα. Η σημασία τους βρίσκεται στο ότι το πλήθος των κύκλων που εκτελείται για την εκτέλεσή τους καθορίζει τη διαμεταγωγή αποθήκευσης του Διαχειριστή ουρών. Ο πίνακας 2.1 παρουσιάζει τις διαθέσιμες εντολές του διαχειριστή ουρών τα ορίσματά τους, τα δεδομένα που επιστρέφουν και το πλήθος των κύκλων που χρειάζονται για την εκτέλεσή τους.

Table 2-1: Πίνακας όλων των εντολών του Διαχειριστή Ουρών

Όνομα	Από	Ορί- σματα	Δεδ. Επιστρ.	Κυ- κλοι	Περιγραφή
Read	CPU	Address	Mem Data	5	Διαβάζει μια λέξη 64 bit από την μνήμη SDRAM
Write	CPU	Address, Data		5	Γράφει μια λέξη 64 bit από στην μνήμη SDRAM
OpenFl	CPU	FlowID, FGID, Hwmark, LWmark		10	Αρχικοποιεί μια ροή κατά την έναρξη μιας σύνδεσης. Δεσμεύει ένα Flow ID και αναθέτει τη ροή σε μία ομάδα ροών. Θέτει τις παραμέτρους ελέγχου ροής για τη σύνδεση (Hwmark, Lwmark).
CloseFl	CPU	FlowID		20	Κλείνει μια ροή κατά τη διάλυση μιας σύνδεσης. Ελευθερώνει το Flow ID και το αφαιρεί από την κυκλική λίστα
ReadCnt	CPU	FlowID	Counter	5	Διαβάζει το μετρητή των αποθηκευμένων κελιών μιας ροής από την εγγραφή της στην μνήμη
Enqueue	Cell Demux	FlowID, Cell		20, 40	Εισάγει ένα εισερχόμενο κελί στην αντίστοιχη ουρά του Flow ID που έχει στην επικεφαλίδα του
Dequeue	Cell Sched	FGID	Address, Cell	20, 40	Εξάγει ένα κελί από την κεφαλή της Ροής που είναι κεφαλή της κυκλικής λίστας της ομάδας ροής FGID. Το στέλνει στον Cell Scheduler μαζί με την διεύθυνση του χώρου που δεσμεύει στην SDRAM.
RdCell	Cell Sched	Address		12	Διαβάζει τα περιεχόμενα ενός χώρου κελιού στην SDRAM χρησιμοποιώντας την διεύθυνση που δίνεται.
Free	Cell Sched	Address		10	Βάζει τον χώρο κελιού της διεύθυνσης που δίνεται στη λίστα ελευθέρων.
ChParam	CPU	FlowID, Hwmark, LWmark		10	Αλλάζει τις παραμέτρους του Ελέγχου ροής της σύνδεσης FlowID σε αυτές που δίνονται.

Η εντολή Enqueue απαιτεί 20 κύκλους για την εκτέλεσή της εκτός της περίπτωσης που η ροή ήταν ανενεργή. Σε αυτή την ειδική περίπτωση η ροή πρέπει να εισαχθεί στην κυκλική λίστα της ομάδας ροής που ανήκει. Οι εγγραφές της επόμενης και προηγούμενης ροής σε αυτή τη λίστα πρέπει να ενημερωθούν προσθέτοντας άλλους 20 κύκλους στο σύνολο των κύκλων εκτέλεσης. Έτσι το πλήθος τους αυξάνεται σε 40.

Η εντολή Dequeue χρειάζεται επίσης 20 κύκλους για την εκτέλεσή της, εκτός της ειδικής περιπτώσεως που η ροή γίνεται ανενεργή (το κελί που εξάχθηκε ήταν το τελευταίο της ροής στη μνήμη SDRAM). Σε αυτή τη περίπτωση η ροή πρέπει να αφαιρεθεί από την κυκλική λίστα της ομάδας ροών που ανήκει. Οι εγγραφές της προηγούμενης και επόμενης ροής πρέπει να ενημερωθούν, προσθέτοντας άλλους 20 κύκλους στο πλήθος των κύκλων εκτέλεσης, που ανεβαίνει στους 40.

2.5 Μορφή της Εγγραφής Ροής

Στο σχήμα 2.4 δίνεται η λεπτομερής περιγραφή των εγγραφών ροών. Υπάρχουν 64 χιλιάδες εγγραφές σαν και αυτή, μία για κάθε μία από τις 64 χιλιάδες ροές. Κάθε εγγραφή αποτελείται από 2 λέξεις των 64 bit. Αυτές είναι 2 λέξεις της μνήμης SDRAM. Οι εγγραφές είναι τοποθετημένες στη μνήμη SDRAM με ευθυγράμμιση 2 λέξεων.

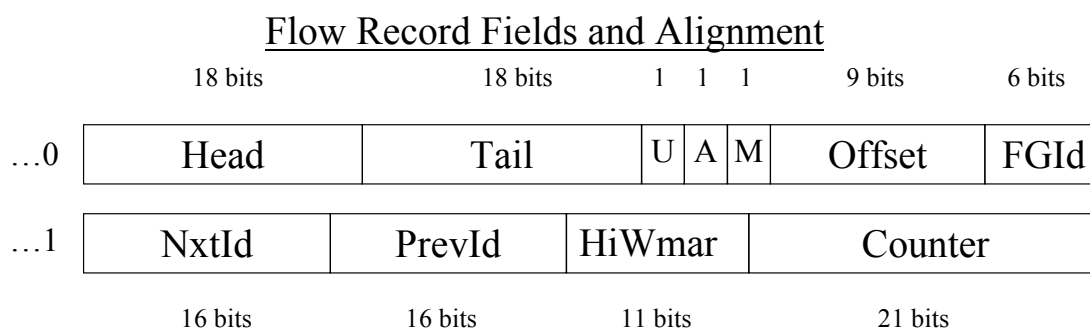


Figure 2-4: Τα πεδία και η ευθυγράμμιση της εγγραφής ροής

Ο πίνακας 2.2 δίνει την περιγραφή κάθε πεδίου της εγγραφής ροής:

Table 2-2 : Πεδία της εγγραφής ροής και περιγραφή τους

Πεδίο	Bits	Περιγραφή
Head	22	Δείκτης Κεφαλής (Head Pointer): Περιέχει τη διεύθυνση της κεφαλής της λίστας της ροής.
Tail	22	Δείκτης Ουράς (Tail Pointer): Περιέχει τη διεύθυνση της ουράς της λίστας της ροής.
Used	1	Χρησιμοποιούμενη Ροή (Used Flow): Όταν είναι 1 αυτή η ροή χρησιμοποιείται από κάποια σύνδεση. Όλα τα πεδία της εγγραφής είναι έγκυρα.

Active	1	Ενεργή Ροή (Active Flow): Όταν είναι 1 αυτή η ροή έχει κάποιο κελί αποθηκευμένο μέσα στην SDRAM.
Mark	1	Μαρκάρισμα (Mark Bit): Όταν είναι 1 το πλήθος των αποθηκευμένων κελιών της ροής στην SDRAM έχει ξεπεράσει το Hi Watermark. Στα εξερχόμενα κελία γίνεται μαρκάρισμα RM, EFCI. Το πεδίο γίνεται 0 όταν ο Counter πέσει κάτω από το Hi Watermark - Off
HiWmH	6	Hi Watermark Most Significant bits: Τα πιο σημαντικά bit του πεδίου Hi Watermark.
HiWmL	12	Hi Watermark Least Significant bits: Τα λιγότερο σημαντικά bit του πεδίου Hi Watermark.
FGId	6	Flow Group ID: Ο αριθμός ταυτότητας της ομάδας ροών που ανήκει η ροή.
Off	5	Offset: Hi Watermark – Off είναι ίσο με το Low Watermark.
NextID	16	Next Flow ID: Το ID της επόμενης ροής στην κυκλική λίστα της ομάδας ροών. Αν Active = 0 αυτό το πεδίο είναι άκυρο.
PrevID	16	Previous Flow ID: Το ID της προηγούμενης ροής στην κυκλική λίστα της ομάδας ροών. Αν Active = 0 αυτό το πεδίο είναι άκυρο.
Counter	20	Το πλήθος των κελιών της ροής μέσα στην μνήμη SDRAM.

Έγινε προσπάθεια κατά την σχεδίαση, η εγγραφή ροής να είναι 2 λέξεις μνήμης σε μέγεθος. Αν ήταν 3 λέξεις θα χαλούσε την ευθυγράμμιση των λέξεων στην μνήμη SDRAM. Αν ήταν 4 λέξεις θα χρειαζόντουσαν 2 επιπλέον κύκλοι για την προσπέλασή του, κάτι που θα έριχνε την διαμεταγωγή του διαχειριστή ουρών. Οι δείκτες σε χώρους κελιών όπως Head and Tail είναι ευθυγραμμισμένοι στο μεγεθός τους (8x64) στην SDRAM μεγέθους 256 MB, οργανωμένη σε 2^{25} λέξεις των 64 bit. Μιας και κάθε χώρος κελιού είναι 8x64, ένας δείκτης σε ευθυγραμμισμένο buffer είναι 22 bit σε μέγεθος.

Δεν υπάρχει πεδίο για το Flow ID κάθε ροής. Αυτό γίνεται γιατί κάθε εγγραφή ροής είναι αποθηκευμένη σε 2 λέξεις η διεύθυνση των οποίων έχει την τιμή του Flow ID. Ένα επιπλέον bit χρησιμοποιείται για το διαχωρισμό των 2 λέξεων της εγγραφής. Με αυτή την οργάνωση βάζουμε τις εγγραφές στην αρχή του χώρου της SDRAM και γλυτώνουμε σε χώρο εγγραφών.

2.6 Μορφή Εγγραφής ομάδας ροών (Flow Group Record)

Όλη η απαραίτητη πληροφορία για τις ομάδες ροών είναι αποθηκευμένη σε μία μνήμη SRAM 64x33. Κάθε μία από τις 64 λέξεις της μνήμης αυτής αποθηκεύει την

εγγραφή ομάδας ροών (Flow Group Record) με τον αριθμό ταυτότητας ως διεύθυνση της. Διατηρεί όλα τα απαραίτητα δεδομένα για την διατήρηση της κυκλικής λίστας της αντίστοιχης ομάδας ροών. Το σχήμα 2.5 παρουσιάζει την μνήμη των ομάδων ροών και την μορφή των αντίστοιχων εγγραφών ομάδων ροών.

	16 bits	16 bits	1
0	Head Flow	Tail Flow	A
1	Head Flow	Tail Flow	A
2	Head Flow	Tail Flow	A
3	Head Flow	Tail Flow	A

62	Head Flow	Tail Flow	A
63	Head Flow	Tail Flow	A

Figure 2-5: Οργάνωση της μνήμης ομάδων ροών και των εγγραφών ομάδων ροών.

Ο πίνακας 2.3 δίνει τη περιγραφή και το μέγεθος των εγγραφών ομάδων ροών.

Table 2-3: Περιγραφή των πεδίων των εγγραφών ομάδων ροών

Πεδίο	Bits	Περιγραφή
Head Flow	16	Head Flow ID: Περιέχει το FlowID της ροής που θα δώσει το επόμενο κελί, όταν μια εντολή dequeue ζητηθεί από την αντίστοιχη ομάδα ροών..
Tail	16	Tail Flow ID: Περιέχει το Flow ID της ροής που εξυπηρετήθηκε κατά την τελευταία εντολή dequeue από την αντιστοιχή ομάδα ροών.
Active	1	Active Flow: Όταν είναι 1 η ομάδα ροών έχει ενεργές ροές, αλλιώς όλη η ομάδα ροών είναι ανενεργή.

Τα περιεχόμενα της μνήμης ομάδας ροών είναι ορατά και στο Cell Scheduler. Ο τελευταίος χρειάζεται να ξέρει ποιές ομάδες ροών είναι ανενεργές ώστε να τις διατηρεί στο χρονοπρογραμματιστικό αλγόριθμο. Η αίτηση για έξοδο κελιού από μία ανενεργή ομάδα ροών θα προκαλούσε λάθος στο σύστημα.

2.7 Μορφή και ευθυγράμμιση κελιού στη μνήμη SDRAM

Κάθε κελί αποθηκεύεται στην μνήμη SDRAM σε χώρους μεγέθους 8x64. Οι πρώτες 7 λέξεις (7x8=56 bytes) αποθηκεύουν το κελί τύπου ABR μαζί με την εσωτερική στο μεταγωγέα επικεφαλίδα (CellBus, Tandem Routing Header για την περίπτωση του μεταγωγέα Δίπολο). Η τελευταία λέξη αποθηκεύει τον δείκτη στο επόμενο κελί της ουράς της ροής (next pointer). Κάθε χώρος είναι ευθυγραμμισμένος στις 8 λέξεις (8 word alignment). Το σχήμα 2.6 παρουσιάζει την μορφή και την ευθυγράμμιση ενός κελιού μέσα στη μνήμη SDRAM.

...000	Cell 0
...001	Cell 1
...010	Cell 2
...011	Cell 3
...100	Cell 4
...101	Cell 5
...110	Cell 6
...111	Next Ptr

Figure 2-6: Μορφή και ευθυγράμμιση κελιού στη μνήμη SDRAM

2.8 Οργάνωση της μνήμης SDRAM

Η μονάδα μνήμης SDRAM DIMM που καλύπτει της ανάγκες μνήμης του Διαχειριστή Ουρών έχει χωρητικότητα μεγέθους 256 MB. Μέσα στην μνήμη αυτή αποθηκεύονται οι 64 χιλιάδες εγγραφές ροών. Ο υπόλοιπος χώρος χωρίζεται σε χώρους κελιών που ανατίθεται δυναμικά στην εισερχόμενη κίνηση για αποθήκευση κελιών. Οι εγγραφές ροών είναι στατικά αναθεμεμένες. Αυτό σημαίνει ότι όλες οι 64 χιλιάδες εγγραφές είναι παρούσες, ακόμη και αν δεν χρησιμοποιούνται από κάποια ροή.

Η CPU μπορεί να αρχικοποιήσει την μνήμη SDRAM με τη χρήση της εντολής Write. Ακόμη υπάρχει μια μηχανή καταστάσεων που μετά από reset μπορεί να αρχικοποιήσει την μνήμη SDRAM. Τα περιεχόμενα της μνήμης SDRAM μετά από την αρχικοποίηση φαίνονται στο σχήμα 2.7.

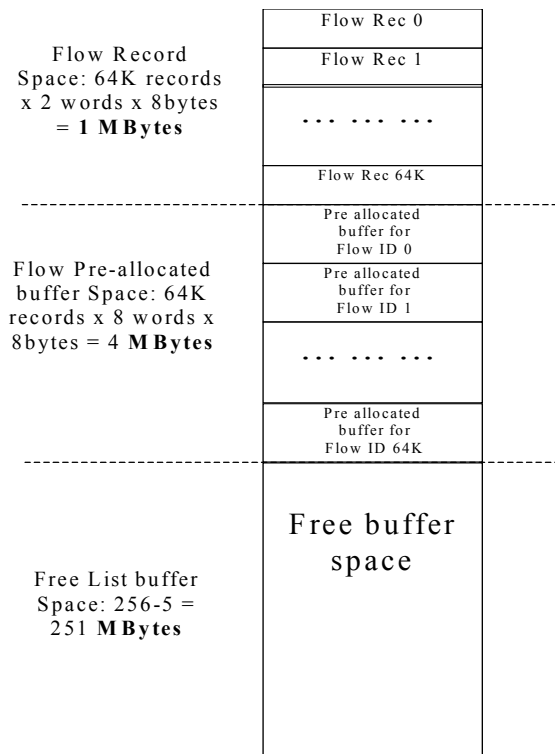


Figure 2-7: Καταμερισμός και οργάνωση χώρου της μνήμης SDRAM

Μιας και η μνήμη είναι $256 \text{ MB} = 2^{28} \text{ bytes}$, οργανωμένα σε λέξεις των 64 bit (8 bytes), το πλήθος των συνολικών λέξεων στην μνήμη είναι 2^{25} . Οι δείκτες σε χώρους κελιών που είναι ευθυγραμμισμένοι στις 8 λέξεις είναι κατά συνέπεια 22 bits. Όπως περιγράφεται στην παράγραφο 2.5 προκειμένου να μην υπάρχει πεδίο για το FlowID μιας εγγραφής, το τελευταίο χρησιμοποιείται σαν δείκτης στη θέση της εγγραφής ροής στην μνήμη. Οι εγγραφές τοποθετούνται στην αρχή της μνήμης. Έτσι η διεύθυνση της πρώτης λέξης της εγγραφής με FlowID 0b1111000111110001 (Το FlowID είναι μεγέθους 16 bit) είναι 0b{00000000,1111000111110001,0}, ενώ η διεύθυνση της δεύτερης λέξης είναι 0b{00000000,1111000111110001,1}. Μιας και ο Διαχειριστής Ουρών υποστηρίζει 2^{16} ροές (64 χιλιάδες), οι εγγραφές ροών δεσμεύουν συνολικά το πρώτο $2^{16} \text{ Flows} * 2 \text{ words/Flow} = 2^{17} \text{ words} = 2^{20} \text{ bytes} = 1 \text{ Mbyte}$ της μνήμης.

Αφού κάθε ουρά ροής πρέπει να έχει ένα κενό χώρο κελιού στο τέλος της (ακόμη και αν δεν χρησιμοποιείται, δεξ παράγραφο 2.11) λόγω της προανάθεσης χώρου, ένας δίνεται σε κάθε μια ουρά κατά την αρχικοποίηση. Αυτοί οι χώροι δεσμεύονται στην μνήμη SDRAM μετά τον χώρο των εγγραφών ροών. Μετά την εκκίνηση του συστήματος αυτοί οι χώροι χρησιμοποιούνται απο εισερχόμενη κίνηση αλλά άλλοι πέρνουν την θέση τους σαν άδειοι χώροι στο τέλος των ουρών. Έτσι το πλήθος των άδειων προανατεθειμένων χώρων κελιών είναι σταθερό και ίσο με 2^{16} (ένα για κάθε υποστηριζόμενη ροή). Η προανάθεση χώρων, δηλαδή, χρησιμοποιεί άλλα $2^{16} \text{ buffers} * 8 \text{ words/buffer} = 2^{19} \text{ words} = 2^{22} \text{ bytes} = 4 \text{ MByte}$.

Τα υπόλοιπα $256 - (4+1) = 251 \text{ MByte}$ χρησιμοποιούνται για αποθήκευση κελιών. Κατά την αρχικοποίηση της μνήμης αυτοί οι χώροι οργανώνονται σε μια μεγάλη λίστα ελευθέρων μονής σύνδεσης (Free list). Αυτή η λίστα χρησιμοποιείται για να παρέχει χώρους σε εντολές enqueue και να δέχεται χώρους μετά από εντολές dequeue και free.

2.9 Μηχανές Καταστάσεων (State Machines)

Όπως περιγράφεται στην παράγραφο 2.1, η υπομονάδα μηχανής καταστάσεων του διαχειριστή ουρών περιέχει τις μηχανές καταστάσεων που εκτελούν τις εντολές που ζητούνται από άλλες μονάδες.

Το σχήμα 2.8 δείχνει την εσωτερική ιεραρχία αυτών των μηχανών καταστάσεων μέσα στην υπομονάδα.

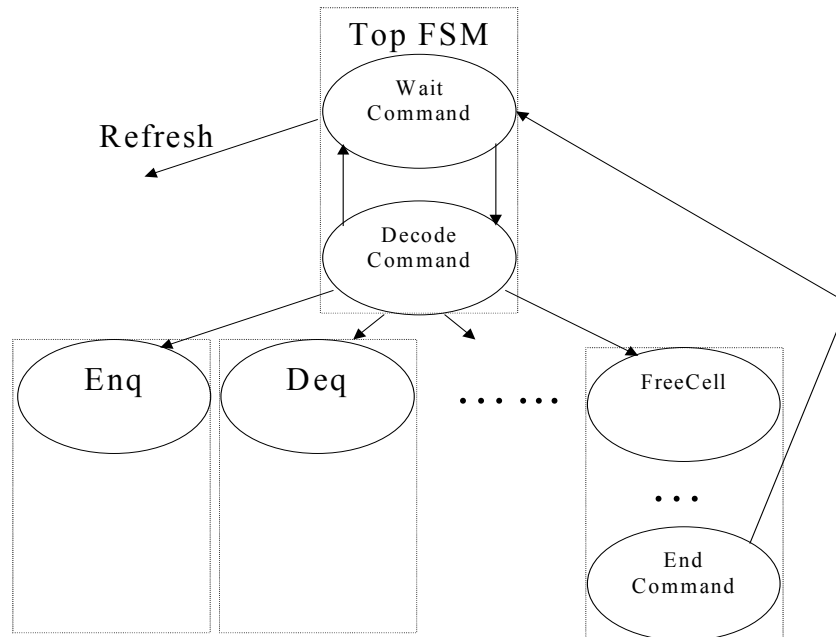


Figure 2-8: Διάγραμμα καταστάσεων των μηχανών καταστάσεων

Αυτές οι μηχανές θέτουν τα σήματα ελέγχου του Datapath του διαχειριστή ουρών. Θέτουν επίσης και λαμβάνουν σήματα ελέγχου από τις διεπαφές του διαχειριστή με τις άλλες μονάδες και ζητούν προσπελάσεις στην SDRAM από τον ελεγχτή SDRAM. Υπάρχει μία μηχανή καταστάσεων για κάθε εντολή του διαχειριστή ουρών. Στην κορυφή της ιεραρχίας βρίσκεται η μηχανή καταστάσεων TOP FSM. Οι λειτουργίες αυτής της μηχανής είναι οι εξής:

- **Αρχιμοποίηση καταχωρητών του Datapath μετά από reset του συστήματος:** Μερικοί καταχωρητές (Freelist Head) πρέπει να αρχικοποιηθούν με μια συγκεκριμένη τιμή μετά από reset. Η μηχανή θέτει τα σήματα ελέγχου για αυτό .
- **Αποδοχή των αιτήσεων εντολών από άλλες υπομονάδες:** Η μηχανή παραμένει σε αδράνη περιμένοντας σήματα αιτήσεων για την εκτέλεση εντολών από τις υπομονάδες που τις γεννούν.
- **Διαιτέυση της εκτέλεσης των εντολών του Διαχειριστή Ουρών:** Υπάρχει η περίπτωση, παραπάνω από μία μονάδες να ζητούν την εκτέλεση μιας εντολής, την ίδια στιγμή. Η μηχανή TOP διαιτεύει πια από όλες τις ταυτόχρονες αιτήσεις θα ικανοποιηθεί με βάση κάποιες προτεραιότητες. Όταν μια εντολή πρόκειται να εκτελεστεί ,η μηχανή θέτει το σήμα ελέγχου που ενεργοποιεί την αντίστοιχη μηχανή καταστάσεων για αυτή την εντολή. Στη συνέχεια περιμένει το σήμα ελέγχου από την μηχανή της εντολή που δηλώνει ότι η εντολή ολοκληρώθηκε. Παράλληλα, επιβεβαιώνει την αποδοχή της εντολής στη μονάδα που τη γέννησε.

- **Αίτηση για φρεσκάρισμα (Refresh) της μνήμης SDRAM:** Οι μνήμες SDRAM πρέπει να φρεσκάζονται περιοδικά. Η μηχανή έχει ένα εσωτερικό μετρητή που όταν φτάσει σε μια τιμή, ζητάει από τον ελεγχτή SDRAM μια εντολή refresh.

Table 2-4: Προτεραιότητες των εντολών του Διαχειριστή ουρών

Προτεραιότητα	Όνομα	Από	Σχόλια
1 (highest)	Refresh	State Machine	Η refresh πρέπει να εκτελείται περιοδικά αλλιώς χάνονται δεδομένα από τη μνήμη.
2	Write, Read, OpenFl, CloseFl, ReadCnt, ChParam	CPU	Οι εντολές της CPU υψηλότερη προτεραιότητα από αυτές που γεννιούνται από τους Cell Demux, Cell Scheduler γιατί είναι σύντομες και διαμορφώνουν και ελέγχουν τη λειτουργία του διαχειριστή ουρών και είναι αναγκαίες για την διόρθωση λαθών.
3	Enqueue,	Cell Demux	Η εντολή enqueue (που τίθεται από τον Cell Demux) έχει ίση προτεραιότητα με τις εντολές που τίθενται από τον Cell Scheduler. Όταν υπάρχει ανταγωνισμός ανάμεσα στην εντολή enqueue και μια εντολή του cell Scheduler, μια εναλλάσσουσα διατιήτευση γίνεται από την TOP μηχανή καταστάσεων. Αν η πιο πρόσφατη εκτέλεση ήταν της εντολής enqueue τότε η εντολή του Cell Scheduler, αλλιώς αν η πιο πρόσφατη ήταν του Scheduler τότε εκτελείται η εντολή enqueue.
	Dequeue, RdCell, Free	Cell Scheduler	

Ο πίνακας 2.4 δίνει τις προτεραιότητες των εντολών που λαμβάνονται υπόψιν όταν υπάρχει ανταγωνισμός που πρέπει να επιλυθεί από την μηχανή καταστάσεων TOP. Η προτεραιότητα των εντολών δίνεται από την μεγαλύτερη προς τη μικρότερη. Εντολές στην ίδια γραμμή έχουν την ίδια προτεραιότητα.

Το σχήμα 2.9 δίνει το απλοποιημένο διάγραμμα καταστάσεων της μηχανής καταστάσεων που εκτελεί τη εντολή enqueue. Αυτή η μηχανή μαζί με αυτή της εντολής dequeue είναι οι πιο σύνθετες και πιο συχνές στην υπομονάδα FSM. Οι καταστάσεις στις οποίες γίνεται μια αίτηση για προσπέλαση στην SDRAM στον ελεγχτή SDRAM, φαίνονται στο σχήμα με μεγάλους κύκλους και περιέχουν το είδος της προσπέλασης. Το μέγιστο πλήθος των καταστάσεων/κύκλων ανέρχεται στους 40. Αυτή είναι η περίπτωση που η ροή ήταν προηγουμένως ανενεργή. Τότε είναι που η ροή θα πρέπει να επαναεισαχθεί στην κυκλική λίστα της ομάδας ροής που ανήκει. Οι εγγραφές της προηγούμενης και επόμενης ροής στη κυκλική λίστα πρέπει να ενημερωθούν και αυτό κοστίζει 20 επιπλέον κύκλους ρολογιού στην εκτέλεση της εντολής. Έτσι το πλήθος των κύκλων ανέρχεται στους 40. Άλλοι 5 κύκλοι

γλυτώνονται όταν υλοποιείται η παράκαμψη λίστας ελευθέρων (δες παράγραφο 2.10). Σε αυτή τη περίπτωση, μια πρόσβαση στη λίστα ελευθέρων αποφεύγεται. Οι προσπελάσεις της μνήμης SDRAM γίνονται με τέτοια σειρά ώστε να αποφεύγονται οι εξερτήσεις δεδομένων (το αποτέλεσμα μιας προσπέλασης να χρησιμοποιείται σαν διεύθυνση της αμέσως επόμενης προσπέλασης), ενώ παράλληλα η μνήμη χρησιμοποιείται συνεχώς.

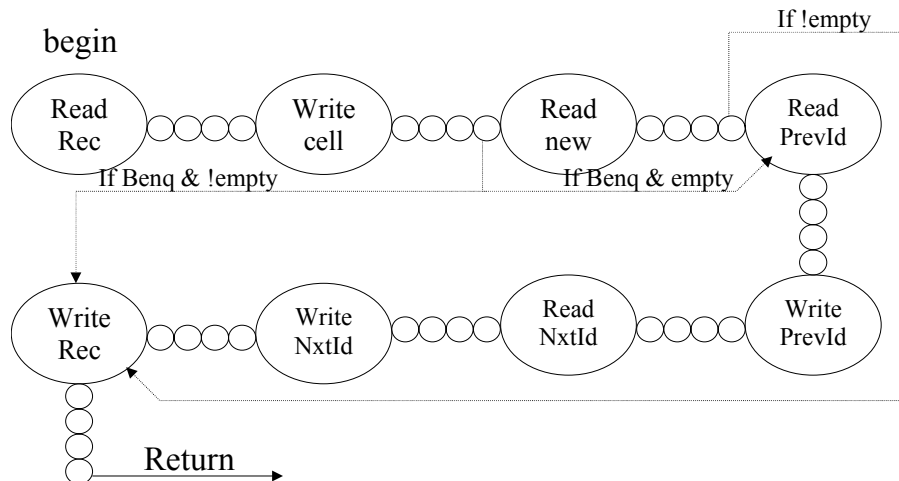


Figure 2-9: Διάγραμμα καταστάσεων της μηχανής καταστάσεων της εντολής Enqueue

2.10 Παράκαμψη λίστας ελευθέρων (Free List Bypassing) [13]

Η παράκαμψη της λίστας ελευθέρων είναι μια τεχνική αποθήκευσης που υλοποιείται από τον διαχειριστή ουρών για να αποφεύγει κάποιες προσβάσεις στην μνήμη SDRAM κατά την διάρκεια μιας enqueue και dequeue εντολής. Έτσι μειώνονται οι απαιτούμενοι κύκλοι για την εκτέλεση και των δύο και έτσι αυξάνεται η διαμεταγωγή του διαχειριστή ουρών.

Όταν μια εντολή enqueue εκτελείται, ένα νέος χώρο κελιού πρέπει να δοθεί στην ουρά που λαμβάνει το κελί. Ο δείκτης κελιών που δείχνεται από την κεφαλή της λίστας ελευθέρων επιλέγεται γι' αυτό το λόγο αλλά αυτό σημαίνει ότι η κεφαλή πρέπει πλέον να δείχνει στο επόμενο χώρο στη λίστα ελευθέρων. Η ανάγνωση αυτού του δείκτη από το χώρο που δίνεται στην enqueue κοστίζει 5 κύκλους μίας και βρίσκεται στην SDRAM.

Ακόμη, κατά την εκτέλεση μιας εντολής dequeue το εξερχόμενο κελί, αν μεταδοθεί σωστά, μπορεί πλέον να αποδεσμεύσει το χώρο του στην μνήμη ώστε να μπει στη λίστα στο χώρο ελευθέρων. Αυτό γίνεται στον διαχειριστή ουρών με την εκτέλεση της εντολής Free που ενημερώνει τον δείκτη κεφαλής της λίστας ελευθέρων με την διεύθυνση του νέου κελιού και την εγγραφή στον next pointer αυτού του χώρου την διεύθυνση του προηγούμενου χώρου που ήταν κεφαλή. Αυτή η εγγραφή κοστίζει 5 κύκλους για κάθε έξοδο κελιού, μιας και γίνεται στη μνήμη SDRAM.

Η παράκαμψη της λίστας ελευθέρων αποφεύγει τους 5 κύκλους που ξοδεύονται στην είσοδο και έξοδο ενός κελιού που περιγράφηκαν πιο πάνω. Αντί οι χώροι που ελευθερώνονται μετά από την έξοδο ενός κελιού να τοποθετούνται στην λίστα ελευθέρων, η διεύθυνσή τους μπαίνει σε μια FIFO εσωτερική στον Cell Scheduler. Όταν μια επόμενη εντολή enqueue ξεκινήσει την εκτέλεσή της, ο νέος χώρος κελιού

που ζητείται, δίνεται από αυτή τη FIFO και έτσι η πρόσβαση στη λίστα ελευθέρων αποφεύγεται.

Συνολικά η παράκαμψη της λίστας ελευθέρων μειώνει το πλήθος των κύκλων μιας εντολής enqueue και dequeue κατά 10 κύκλους ρολογιού. Μιάς και οι δύο εντολές χρειάζονται 20 ή 40 κύκλους για την εκτέλεση τους, αυτή η τεχνική βελτιώνει την επίδοση του διαχειριστή ουρών κατά $10/(20+20) = 1/4 = 25\%$ ή $10/(40+40) = 1/8 = 12,5\%$.

Το σχήμα 2.10 παρουσιάζει την παράκαμψη της λίστας ελευθέρων.

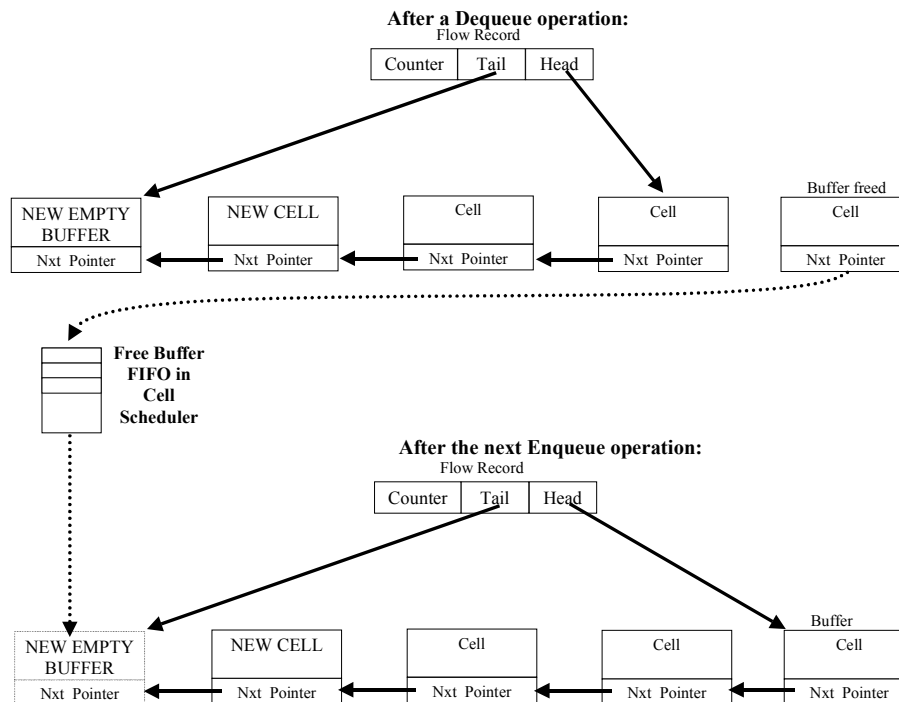


Figure 2-10: Υλοποίηση της παράκαμψης λίστας ελευθέρων στον διαχειριστή ουρών

2.11 Προανάθεση Χώρου κελιών (Cell Buffer Pre-allocation) [10]

Η προανάθεση χώρου κελιών είναι μια τεχνική που χρησιμοποιείται από το διαχειριστή ουρών που με κόστος ένας αχρησιμοποίητο χώρο κελιού για κάθε μία υποστηριζόμενη ροή, μειώνει το πλήθος των κύκλων μιας εντολής enqueue κατά 5.

Όπως περιγράφηκε στις προηγούμενες παραγράφους, κάθε μία από τις 64 χιλιάδες ουρές του διαχειριστή ουρών έχει πάντα στο τέλος της ένα άδειο στοιχείο. Αυτό ισχύει ακόμη και για τις άδεις (ανεργές) ουρές. Σε αυτή τη περίπτωση τα πεδία Head και Tail της αντίστοιχης εγγραφής δείχνουν και τα δύο στον άδειο χώρο κελιού. Αυτός ο χώρος ονομάζεται προαναθεμένος χώρος. Αν ο χώρος δεν ήταν προαναθεμένος και ο δείκτης Tail έδειχνε στο τελευταίο χρησιμοποιούμενο χώρο, κατά την διάρκεια μιας enqueue, το νέο κελί θα έπρεπε να γραφτεί σε ένα νέο άδειο χώρο κελιού, από τη λίστα ελευθέρων (ή από το μηχανισμό παράκαμψης ελευθέρων) και ο δείκτης του προηγούμενου χώρου θα έπρεπε να ενημερωθεί με τη διεύθυνση του νέου χώρου. Μιας και με αυτό το μηχανισμό, γίνονται προσβάσεις σε δύο

ξεχωριστους χώρους, δύο ξεχωριστές γραμμές της SDRAM θα πρέπει να τροποποιηθούν.

Αντίθετα με την προανάθεση χώρου κελιών, το κελί γράφεται στο προαναθεμένο άδειο χώρο κελιού, που είναι τελευταίο στην ουρά και ο next pointer του χώρου γράφεται με την διεύθυνση του νέου άδειου χώρου που δεσμεύεται από τη λίστα ελευθέρων (ή από το μηχανισμό παράκαμψης ελευθέρων). Μόνο ένας χώρος κελιου προσπελαύνεται με αυτό το μηχανισμό γλυτώνοντας 5 κύκλους σε σχέση με τον προηγούμενο μηχανισμό.

Συνολικά η προανάθεση χώρου κελιών μειώνει το πλήθος των κύκλων μιας εντολής enqueue και dequeue κατά 5 κύκλους ρολογιού. Μιάς και οι δύο εντολές χρειάζονται 20 ή 40 κύκλους για την εκτέλεση τους, αυτή η τεχνική βελτιώνει την επίδοση του διαχειριστή ουρών κατά $5/(20+20) = 1/8 = 12,5\%$ ή $5/(40+40) = 1/16 = 6,25\%$.

Το κόστος σε μνήμη της τεχνικής είναι 4 Mbytes που είναι $4/256 = 1/64 = 1,56\%$ της συνολικής μνήμης.

2.12 Θέματα χρονισμού

Σε αυτό το υποκεφάλαιο εξετάζονται οι ικανότητες διαμεταγωγής του διαχειριστή ουρών σε συνάρτηση με τις απαιτήσεις διαμεταγωγεί εισόδου εξόδου της κάρτα εξυπηρετητή ABR. Οι παράμετροι χρονισμού που υποθέτονται αποδεικνύουν ότι η εφαρμογή των τεχνικών Free list bypassing μαζί με την Cell Buffer pre-allocation είναι απαραίτητες, αν ο διαχειριστής ουρών θέλει να επιτύχει στην ικανοποίηση των απαιτήσεων της κάρτας χρησιμοποιώντας μια μοναδική μνήμη SDRAM για όλες τις ανάγκες του.

2.12.1 Το ρολόι UTOPIA σε σχέση με το ρολόι του Διαχειριστή ουρών (ABRSU).

Προκειμένου η αρχιτεκτονική αποθήκευσης που υλοποιείται στον διαχειριστή ουρών να μπορεί να εξισορροπήσει την εισερχόμενη και εξερχόμενη κίνηση του μεταγωγέα, πρέπει να μπορεί να εισάγει στη και να εξάγει απο την μνήμη ένα κελί στο χρόνο άφιξης ενός κελιού (cell time). Ο χρόνος άφιξης κελιού είναι ο χρόνος που απαιτείται για την πλήρη είσοδο ενός κελιού στις διεπαφές UTOPIA. Μιαςκαι οι εντολές enqueue και dequeue δεν μπορούν να εκτελεστούν παράλληλα ο διαχειριστής ουρών πρέπει να χρησιμοποιήσει επιτάχυνση του ρολογιού του σε σχέση με το ρολόι της διεπαφής UTOPIA.

Ένα κελί χρειάζεται 28 κύκλους ρολογιού UTOPIA για να εισέλθει στην ABRSU (διεπαφή πλάτους 16 bit). Οι εντολές enqueue και dequeue χρειάζονται:

80 (Worst case) ή 40 (Normal Case) κύκλους ρολογιού του Διαχειριστή ουρών.

Αυτό σημαίνει ότι η επιτάχυνση του ρολογιού για τις 2 περιπτώσεις είναι:

Worst case:

1 Cell arrival time = 1 Enq time + 1 Deq time=>

$Tclk_utopia * 28 utopia_cycles = Tclk_qm * (40+40) qm_cycles =>$

$Tclk_utopia / Tclk_qm = 80 / 28 = 2,8 =>$

Speed_up_worst = 2,8

Normal case:

$1 \text{ Cell arrival time} = 1 \text{ Enq time} + 1 \text{ Deq time} \Rightarrow$
 $\text{Tclk_utopia} * 28 \text{ utopia_cycles} = \text{Tclk_qm} * (20+20) \text{ qm_cycles} \Rightarrow$
 $\text{Tclk_utopia} / \text{Tclk_qm} = 40 / 28 = 1,4 \Rightarrow$
 $\text{Speed_up_normal} = 1,4$

Η FPGA που χρησιμοποιείται από τη σχεδίαση μπορεί να πετύχει συχνότητα ρολογιού μέχρι και 50 MHz. Αφού τα ρολόγια των Cubit του μεταγωγέα ΔΠΟΛΟ έχουν συχνότητα 25 MHz, η μέγιστη επιτάχυνση στο ρολόι του διαχειριστή ουρών που μπορεί να επιτευχθεί είναι 2. Σε αυτή την περίπτωση:

$$\text{QM_Clk} = 2 * \text{UTOPIA_Clk}.$$

Μιας και :

$$\text{Speed_up_worst} > \text{Speed_up_sel} = 2 > \text{Speed_up_normal}$$

ο διαχειριστής ουρών είναι η τροχοπέδη του συστήματος όταν οι εντολές enqueue και dequeue χρειάζονται το μέγιστο αριθμό κύκλων εκτέλεσης. Οι διεπαφές UTOPIA είναι η τροχοπέδη στην περίπτωση που οι εντολές enqueue και dequeue χρειάζονται τον ελάχιστο αριθμό κύκλων εκτέλεσης.

2.12.2 Αποτελέσματα σύνθεσης, συνεισφορά της παράκαμψη λίστας ελευθέρων και της προανάθεσης χώρων κελιού.

Μετά από τη σύνθεση των αρχείων της Verilog που περιγράφουν τις μονάδες της ABRSU με το εργαλείο σύνθεσης MaxPlusII, πήραμε τα παρακάτω αποτελέσματα:

- **Ρολόι ABRSU (Διαχειριστή ουρών) στα 35 MHz.** Αυτό το ρολόι αποδίδει συνδιασμένη εισερχόμενη και εξερχόμενη διαμεταγωγή 400Mbps για εντολές enqueue και dequeue 40 κύκλων και 800 Mbps για εντολές enqueue και dequeue 20 κύκλων.
- **Χρησιμοποίηση FPGA SRAM στο 95%**
- **Χρησιμοποίηση λογικών πυλών FPGA στο 55%**

Θεωρώντας 35 MHz συχνότητα ρολογιού του διαχειριστή ουρών η επιτάχυνση είναι 1,4 αντί 2 που θεωρήθηκε στην παράγραφο 2.12.1

Προκειμένου να μην είναι ο διαχειριστής ουρών τροχοπέδη του συστήματος πρέπει:

$1 \text{ Cell arrival time} \geq 1 \text{ Enq time} + 1 \text{ Deq time} \Rightarrow$
 $\text{Tclk_utopia} * 28 \text{ utopia_cycles} \geq \text{Tclk_qm} * (20+20) \text{ qm_cycles} \Rightarrow$
 $\text{Tclk_utopia} / \text{Tclk_qm} \geq 50 / 28 = 1,8 \Rightarrow$
 $40 \text{ ns} / 28\text{ns} \geq 1,4 \Rightarrow$
 $1,4 \geq 1,4$

Αυτό σημαίνει ότι η επιτευχθείσα επιτάχυνση μόλις που ικανοποιεί τις ανάγκες διαμεταγωγής των διεπαφών UTOPIA. Αν η παράκαμψη λίστας ελευθέρων και η προανάθεση χώρων κελιών δεν εφαρμοζόταν τότε οι κύκλοι ρολογιού που θα χρειαζόντουσαν για την κανονική περίπτωση της εκτέλεσης των εντολών enqueue και dequeue θα ήταν $20+20 + 10$ κύκλους που αποφεύγουμε με την παράκαμψη λίστας ελευθέρων + 5 κύκλους που αποφεύγουμε με την προανάθεση χώρων κελιών = 55

κύκλοι. Σε αυτή την περίπτωση η συχνότητα ρολογίου που θα ήταν απαραίτητη είναι περίπου 50 MHz.

Η βελτίωση της συνδιασμένης επίδοσης του συστήματος χάρη στις δύο αυτές τεχνικές είναι $15 \text{ κύκλοι} / 55 \text{ κύκλοι} = 27 \%$.

3 Συμπεράσματα και μελλοντικές επεκτάσεις

Σε αυτή την εργασία σχεδιάσαμε και υλοποιήσαμε την αρχιτεκτονική ενός Ανά-Ροή διαχειριστή ουρών με σκοπό την αποθήκευση της κίνησης τύπου ABR σε ένα μεταγωγέα ATM σε κατάσταση συμφόρησης του. Ο διαχειριστής ουρών υλοποιήθηκε σε μια μεγάλη FPGA που τοποθετήθηκε σε μία από τις κάρτες του μεταγωγέα. Η χρήση της FPGA επέτρεψε το εκτενές τεστάρισμα του συστήματος κατά την ανάπτυξή του και μας έδωσε την ικανότητα να επιβεβαιώσουμε τις υποθέσεις μας για την εφικτή ταχύτητα του συστήματος, στα αρχικά ακόμα στάδια της σχεδίασης.

Χρησιμοποιήσαμε μια μοναδική μνήμη SDRAM DIMM για την αποθήκευση των κελιών και των δεικτών των ουρών. Αυτό μείωσε την χρήση pin και συρμάτων της κάρτα, αποδίδοντας ένα σύστημα χαμηλού κόστους. Ο προσεκτικός προγραμματισμός των προσπελάσεων στη μνήμη SDRAM από το διαχειριστή ουρών απέδειξε ότι η προσέγγιση της μοναδικής μνήμης είναι εφικτή.

Παρόλο που η δυναμική παραχώρηση ανέβασε το πλήθος των προσβάσεων για κάθε εντολή enqueue και dequeue μειώνοντας την διαμεταγωγή της αποθήκευσης, μας επέτρεψε κατά το τεστάρισμα να χρησιμοποιήσουμε το διαχειριστή ουρών για αποθήκευση χιλιάδων κελιών σε μία ουρά ροής και παράλληλα να διατηρήσουμε την ικανότητα να χειριζόμαστε 64 χιλιάδες ροές.

Οι διεπαφές του διαχειριστή ουρών με την εξωτερική CPU μας επέτρεψε να εκσφαλματώσουμε το σύστημα αποδοτικά και να εισάγουμε κίνηση που επιβεβαιώνει την ορθότητα της φυσικής σύνδεσης της FPGA με την μνήμη SDRAM DIMM και τα ολοκληρωμένα CubitPro.

Η χρήση των τεχνικών παράκαμψης λίστας ελευθέρων και προανάθεσης χώρου ελευθέρων αποδείχτηκε απαραίτητη για την επίτευξη του στόχου της διαμεταγωγής αποθήκευσης κοντά στο 1 Mbps από το σύστημά μας. Η βελτίωση της τάξης του 26% στην επίδοση του συστήματος που επιφέραν με την εφαρμογή του απόσβεσε την απώλεια του στόχου των 50 MHz στην αρχή της σχεδίασης του συστήματος. Έτσι η μέγιστη εισερχόμενη και εξερχόμενη διαμεταγωγή που επιτεύχθηκε ανήλθε στα 800 Mbps που είναι αρκετή για τις ανάγκες αποθήκευσης κίνησης τύπου ABR ενός μεταγωγέα της τάξης Gbps.

Η προσθήκη διαφόρων άλλων χαρακτηριστικών μεταγωγέων, είναι ένα ενδιαφέρον θέμα για μελλοντική απασχόληση πάνω στο διαχειριστή ουρών. Για παράδειγμα, η μεγένθυση της εγγραφής ροής από 2 λέξεις των 64 bit σε 4 λέξεις θα μπορούσε να επιτρέψει την προσθήκη επιπλέον πεδίων για κάθε μία από τις 64 χιλιάδες υποστηριζόμενες ροές. Ένα πεδίο New ID θα μπορούσε να προστεθεί που να ήταν προσβάσιμο από την CPU. Αυτό το πεδίο θα αντικαθιστούσε το Header ID των εισερχόμενων κελιών μιας ροής. Έτσι ο διαχειριστής ουρών θα μπορούσε να παρέχει και μετάφραση VP/VC (VP/VC translation) παράλληλα με την ανά-ροή αποθήκευση. Ένα άλλο πεδίο που θα μπορούσε να προστεθεί είναι ένα πεδίο Explicit Rate. Η CPU θα υπολόγιζε το κατάλληλο ρητό ρυθμό της σύνδεσης και θα το χρησιμοποιούσε για να αλλάξει το Explicit Rate μέσα στα κελιά RM της σύνδεσης, αν η προυπάρχουσα τιμή είναι μεγαλύτερη από αυτή που θέλει να επιβάλει η CPU. Με αυτό το τρόπο θα υποστηριζόταν ο έλεγχος ροής τύπου RM explicit rate, με το λογισμικό να αναλαμβάνει τον υπολογισμό των ρυθμών και το υλικό να ενημερώνει τα πεδία των κελιών RM.

Άλλα ενδιαφέροντα χαρακτηριστικά που θα μπορούσαν να υποστηριχθούν είναι μηχανισμοί πετάγματος κελιών των ροών με κακή συμπεριφορά.

ΠΑΡΑΡΤΗΜΑΤΑ – Αγγλική μετάφραση

4 Introduction

4.1 Motivation

The introduction of bandwidth hungry applications both in corporate and client communications is one of the most consistent trends in the networking world. Multimedia applications that ride the Moore's Law on one hand, produce excessive amounts of voice and video data to be relayed through the network, while booming business to business networking contributes small but frequent chunks of data between corporate headquarters. These network requirements are heavy burden to the current network infrastructure that relies mostly on IP protocol and wire cables. While the second is the subject of the last mile problem and requires gradual but enormous investment of funds in a global scale with the introduction of fiber optics, the first will still pose problems in efficient networking because IP is incapable of differentiating among diverse services required by network users. ATM protocol should gradually replace or merge with IP applications along the way to high speed, fiber optic networking.

ATM offers generic Quality of Service (QoS) guaranties to existing networking by differentiating network traffic into several types depending on the requirements of the application that produces it and provides special handling and billing for each. These requirements are:

- Bandwidth – The rate at which, the network must carry an application's traffic.
- Latency – The delay that the application can withstand in the delivery of its data
- Jitter – The variation in latency
- Loss – The percentage of acceptable loss of data.

Each of the family of traffic types supported by ATM networks places most of its interest over some of these requirements and less to the remaining. These types are:

- Constant bit rate (CBR) - This type is used for emulating circuit switching. The bandwidth is constant with time. CBR applications are quite sensitive to jitter but not no much to data loss. Examples of applications that can use CBR are telephone traffic, videoconferencing, and television.
- Variable bit rate–non-real time (VBR–NRT) - This type allows users to send traffic at a rate that varies with time depending on the availability of user information. Statistical multiplexing is provided to make optimum use of network resources. Multimedia e-mail is an example of VBR–NRT.
- Variable bit rate–real time (VBR–RT) - This type is similar to VBR–NRT but is designed for applications that are sensitive to cell-delay variation. Examples for real-time VBR are voice with speech activity detection (SAD) and interactive compressed video.
- Available bit rate (ABR) - This type of ATM services provides rate-based flow control and is aimed at data traffic such as file transfer and e-mail. Although the standard does not require the cell transfer delay and cell-loss ratio to be guaranteed or minimized, it is desirable for switches to minimize delay and loss as much as possible. Depending upon the state of congestion in the network, the

source is required to control its rate. The users are allowed to declare a minimum cell rate, which is guaranteed to the connection by the network.

- Unspecified bit rate (UBR) - This type is the catch-all-other class and is widely used today for TCP/IP.

The last two types of traffic are the most opportunistic of them all since they tend to use the leftovers of the network resources reserved by the others whenever available, thus increasing network efficiency. Still they carry the burden of the most traditional of network applications such as E-mail, FTP, HTML that is the base of contemporary networking. Although they pose some bandwidth requirements, they step aside in times of network congestion. This is translated to extensive buffering (or queueing) of these types of traffic at the network nodes they reside when the latter suffers from congestion.

There are a number of architectures proposed throughout the evolution of switches and routers that address the problem of effective traffic queueing. Large buffers alone can not support the variety of traffic types and the problems of delays and HoL (head of line) blocking that arise. It seems that providing separate queueing for each flow of traffic (per flow queueing) gives the ability to the switch/router to service more accurately the requirements of each one [5], [10], [11], [14], while the provision of resizable buffer space (dynamic queueing) for each flow takes advantage of unused, by dormant flows, buffer space [5]. Switch/router architectures with these characteristics can accommodate high speed network traffic while using off-the-shelf, inexpensive memory modules of Dynamic RAM dropping the memory cost to the minimum of a standard PC.

Dynamic RAM memory modules in all their forms have a wide range of usability in contemporary computing and especially SDRAM which is used almost in any computing device that needs large cheap buffer space, with moderate to high throughput. Best performance is achieved when large chunks of data are moved, making these memories suitable for networking that is now well in its packet (~64 bytes) switching era. Their wide applicability, simple interface, and industry-wide standardization drops their cost substantially, fuelling the addition of more ports to the networking device [1,ch9], [12].

4.2 This thesis and the DIPOLO Switch

In this thesis, we describe the architecture of a Queue Manager IP that supports the features of per flow, dynamic queueing of ABR traffic type of ATM networks. This IP was designed for the purposes of a 1Gps ATM Switch called Dipolo.

It is used to accommodate the queueing of a maximum of 64K flows for a centralized ABR Server Card. The Card uses a large FPGA to host the IP and a single SDRAM DIMM module to store cells and queue pointers. There is also a CPU interface that programs the Queue Manager with parameters of each ABR flow so that the IP can support ABR flow control features such as RM marking and EFCI. We also discuss the general architecture of the DIPOLO Switch and the features of the IP architecture that were used to increase Memory utilization, such as free buffer preallocation [10] and free list bypassing [13]. The IP can also support traffic differentiation, by organizing flows into Flow groups according to service needs and output port, making it flexible even to support CBR or VBR traffic when special care

is taken by the scheduling hardware. Special care was also given to the interfaces of the IP with other IPs such as the Scheduler and the CPU so that they are simple and effective. The single SDRAM interface did not allow for any parallel accesses. Instead Enqueue and Dequeue commands were implemented as a whole instead of smaller atomic commands to increase Memory utilization. A clock speed of 35 MHz was achieved for the FPGA implementation that is translated to a maximum 800 Mbps of combined incoming and outgoing throughput. In ASIC implementations where clock speeds of 133 MHz are feasible this throughput could rise substantially making this architecture suitable for use as part of a Networking Chip [7], [10]. A total of 2,5K of FPGA Logic Elements was used as well as only 2K bits of on chip (FPGA) SRAM.

DIPOLO Switch design was a joined project by the University of Crete, Greece, the National Technical University of Athens, the Foundation of Research and Technology of Crete and Intracom Company. The purpose of this project is the design and manufacturing of ATM switch of 1 Gbps of throughput, for the provision of broadband networking to domestic VDSL users. The main system requirements were, the provision of CBR, VBR, ABR types of ATM traffic. Other factors that were also taken into account were low system cost, use of commercial chips and the optimum division of work between the conglomerates.

Under the scope of DIPOLO design, the writer had the chance to work on the physical design of the ABR Server Card. This card is described in section 5.5. More precisely, he worked on the definition of its organization, its description in concept level. He also defined the internal organization of the ABRSU unit (FPGA) of the card. This unit is described in chapter 6. He designed the Cell Demux (see section 6.3), the Queue Manager, the main issue of this thesis that is described in chapter 7, the interface of the ABRSU with the card's cpu (see section 6.2) and the interface with the SDRAM memory unit. He was also in charge of top level synthesis and verification. He designed tests and demos that proved the feasibility of the design.

4.3 Switch/Router Generations and Queueing Architectures

In this section, a brief description of the switch/routers generations that evolved over the years and the queueing architectures that were introduced with each one of them is given.

4.3.1 First Generation Switches/Routers

The evolution of networking devices can be roughly separated into three generations mainly by the hardware used and the level of integration. The first generation of switches were devices that resembled to a general-purpose computer and consisted of line cards interconnected by a I/O bus. There was also a CPU that hosted all the routing software that decided where to forward the packet, a Main Memory, and an optional DMA module to release the CPU of the burden of moving data packets between line cards and Main Memory. A data packet would enter the device through the ingress side of the Line Card by use of an Analog to Digital Converter. Then the CPU would extract the header of the packet in order to route the packet over the Bus to the appropriate queue in the Main Memory. Later the CPU would schedule the packet to be routed over the BUS again, to the Line Card, through a Digital to Analog

converter and off to its destination. Figure 4-1 depicts the configuration of first generation devices.

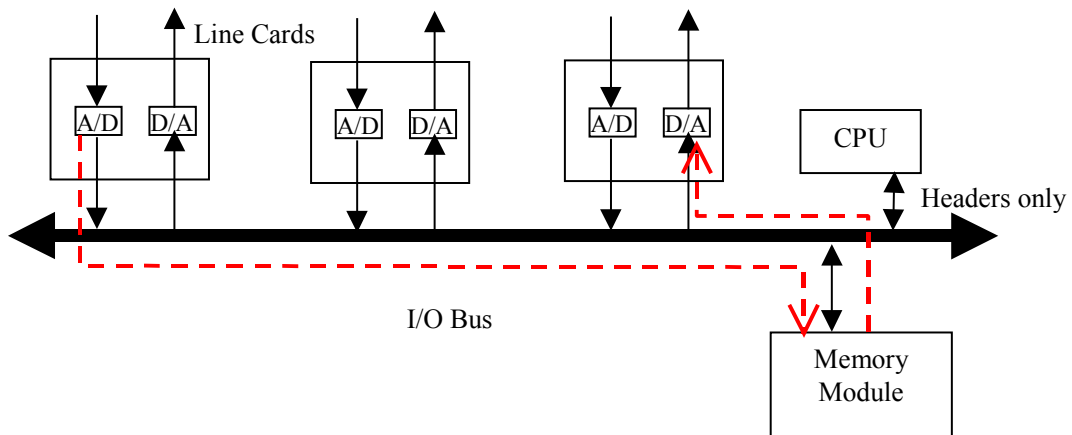


Figure 4-1: First Generation Switch Routers

There are three bottlenecks in this architecture:

- CPU Power
- Memory Throughput
- I/O Bus bandwidth

As the line speed rates increased the CPU would perform poorly, since it implemented both routing and scheduling of packets for all line cards. As for Main Memory and I/O Bus, both would scale badly with the introduction of additional cards especially for the bus that relayed each packet twice through its course through the device. Still, these devices were sufficient for low speed rates and the poor data production of applications, during the first years of the Internet era.

4.3.2 Second Generation Switches/Routers

The second generation of switches eliminated the CPU and Memory bottlenecks of the first by introducing redundancy to both of them. As seen in Figure 4-2, each line card now owns a separate memory module and a small CPU. The main memory is, now, not necessary. The local CPU implements routing and scheduling of packets, while storing takes place in the local memory. Input Queueing or Output Queueing or both can be implemented. The sole purpose of the central CPU is to arbitrate the usage of the Bus, the exchange of routing information between the local cards and the programming and maintenance of the whole system. Due to the redundancy, the only bottleneck of these devices was I/O Bus bandwidth that failed to scale along with the high speed line cards and the port count.

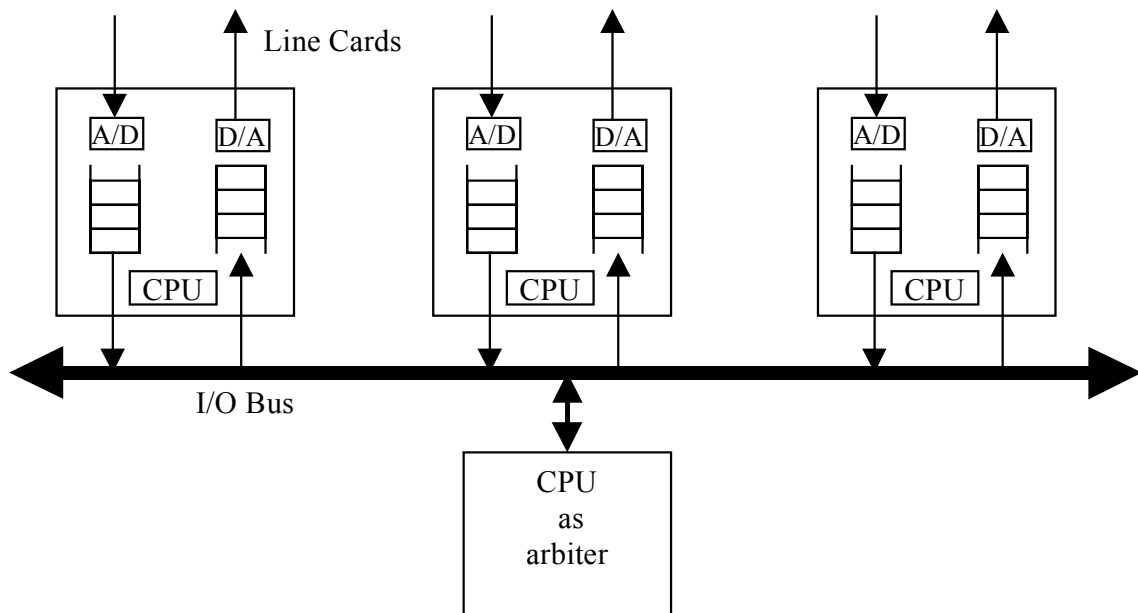


Figure 4-2: Second Generation Switch/Router

4.3.3 Third Generation of Switches/Routers

This generation introduced Switching fabrics to replace the I/O bus as the medium to relay packets between cards. Buffering and routing of data packets is performed inside the line cards while specialized hardware is provided to give to the line cards access to the fabric. Switching fabric can accept multiple simultaneous transfers of packet with a maximum of N transactions when N Line cards are connected to the fabric. The current trend for network devices, that rides the ASIC large-scale integration, is SoC (System on Chip) architectures. Except analog parts, all the line hardware (buffers, routing) for all ports are stored inside the same chip along with a crossbar (the most effective but less scalable of switching fabric), a scheduling unit and a CPU. Such chips can accommodate up to 32 input/output ports and are sufficient for a low-end switch/router. They can also be used as a building block of a large high-end switch/router. In the latter case, they are organized in Switching Fabric topologies such as Banyan, Benes, and Batcher-Banyan networks [1,chapter 8]. All are depicted in figure 4-3.

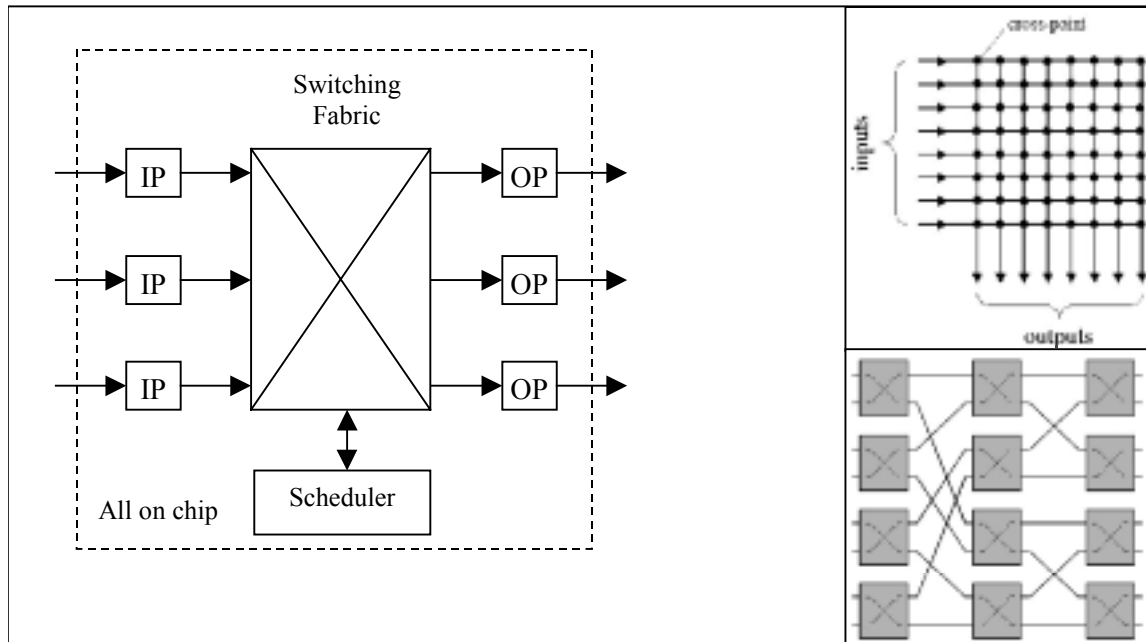


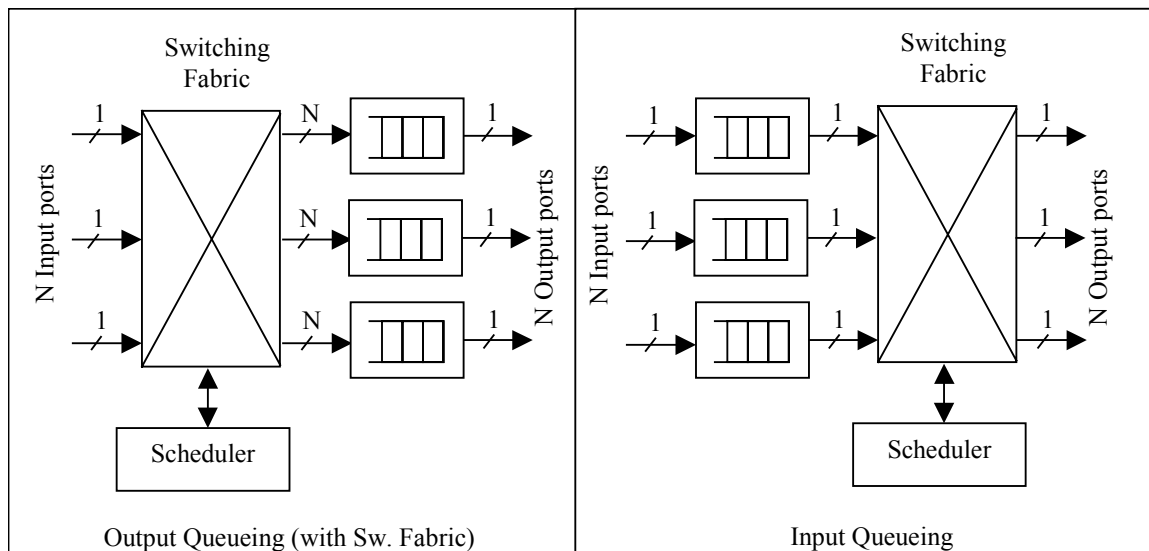
Figure 4-3: Left: Third generation Switch/Router, Top-Right: A crossbar, Bottom-Right: An 8x8 Banyan Fabric made of small 2x2 Switch blocks.

4.4 Queueing Architectures in general

In this section, we describe the main queueing architectures that are implemented in the spectrum of networking device generations described in the previous section. Advantages and disadvantages are also given. Finally per flow queueing is described.

4.4.1 Output Queueing

There are two queueing architecture families: Input queueing and Output queueing architectures. The output queueing families place the buffering memory that stores the queues near the outputs of the device as shown in figure. The interconnecting medium can be either a shared medium (like an I/O Bus) or a switching fabric (like a cross bar). One single module (a shared buffer architecture) or separate ones, one for each output port can serve all the link output queues. The first case with a shared medium is the first generation of switches/routers. Both of them can accept all the available throughput but place heavy requirements on the rate of interconnecting medium and the memories. In the case of N input ports and N output ports with a dedicated memory module for each output port the module must have a throughput of $N+1$, when input throughput for each port is 1. In the case of a single memory module for all output ports, the module must provide a throughput of $N \times (N+1)$. Both configurations, shown in figure 4-4, are not scalable so output queueing is generally not used.



**Figure 4-4: Left: Output Queueing with a Switching Fabric and multiple buffers
Right: Input Queueing with a Switching Fabric.**

4.4.2 Input queueing

When the buffering memory is placed at the inlet side of each port of the switch we describe this scheme as input queueing. Each time a packet arrives it is placed in a queue at this memory and when it gets to the head of the queue it waits for the scheduler to forward it to the output port of destination. The architectures that fall into input queueing are more scalable than the ones that fall into output queueing since their memories must provide twice the input port throughput. Still, they fail to accept all the input rate, due to Head of Line blocking at each input port that decreases by one third of the optimum the switch throughput. Head of Line blocking can be seen in Figure 4-5. When two inlet queues have at their head position a packet destined to the same output port, the fabric can accept only one of them. The other queue remains idle, although there are other packets behind the head that are destined to other output ports and could be served at the same time. To overcome HoL blocking an alternate scheme is often used depicted in figure 4-5. Each packet that arrives at an input port is stored in a separate queue, according to the output port of destination. This scheme is called Advanced Input Queueing or Virtual Output Queueing and theoretically can achieve 100% utilization but requires fast scheduling hardware to find the optimum schedule from inputs to outputs [2].

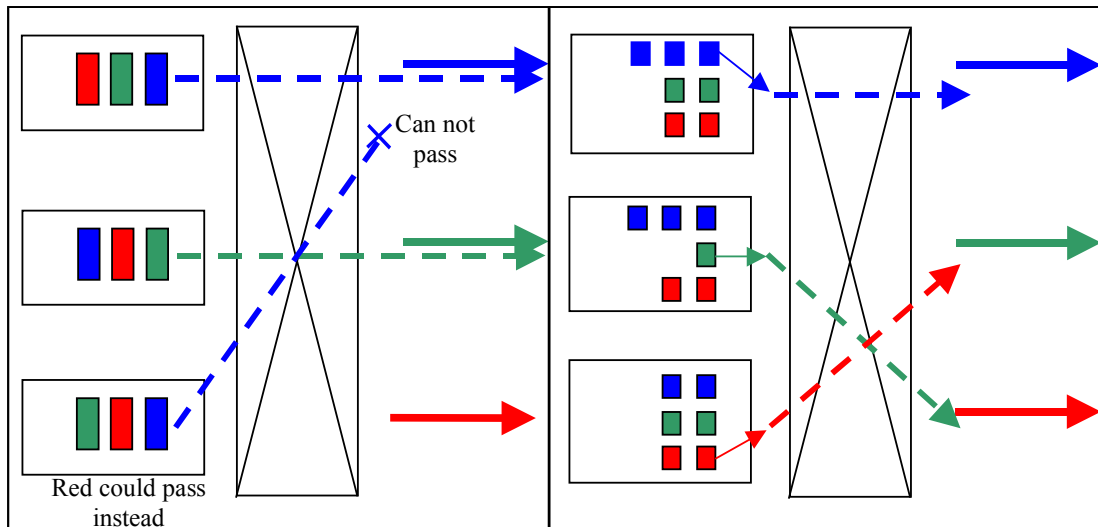


Figure 4-5: Left : Head of Line Blocking Right: Advanced Input Queuing

4.4.3 Variations

Input queueing and Output Queueing can also be used in combination for better performance. In that scheme usually some internal speed-up is used in the switching fabric [3], [4], [5], [6]. The fabric inputs are able to receive packets at a higher rate than the port rate. This resembles to normal Input queueing working at relatively low rates. Additionally, the fabric transmits to the output queues at a higher rate than the output line rate in order to accommodate for the accumulation of packets when most of the traffic at a given time is destined to the same output. Figure 4-6 depicts Internal Speed-Up. Another variation, shown in the figure, is Cross-point or Distributed queueing. It achieves top performance like output queueing by storing packets at the cross-points of a crossbar. This scheme though, requires extensive ($N \times N$) usage of buffers that is costly either on chip or off-chip and is not scalable for large number of ports. Recent advances in embedded DRAM technology (DRAM memory on the same silicon chip as common logic gates) could promote such schemes.

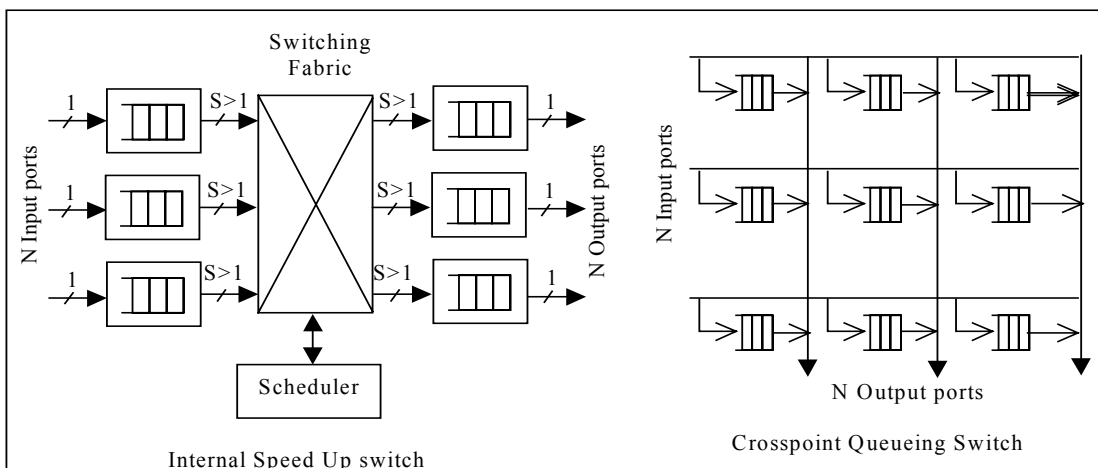


Figure 4-6: Left: Internal Speed Up Switch, Right: Crosspoint Queueing Switch

4.4.4 Per-Flow Queueing Vs Single FIFO Queueing in Output Queueing

In the theory of Queueing there is also the matter of buffer management to be addressed. The buffer size is a limited resource and is heavily contended. The heavy contention over a period of time may lead to congestion, possible loss and extensive delay of packets. Therefore in order to provide Quality of Service, a switch must take care of the buffering and scheduling techniques it uses.

In Output Queueing with a dedicated memory module at each Output, or a shared buffer, the arriving packets could be stored in a single FIFO, or in dedicated queue for each flow of Traffic. In ATM networks where a flow is a VP/VC the latter scheme is called Per-VC queueing.

In the Single queueing approach, a centralized memory is completely shared by a single queue, where all the packets from different sources or input ports enter. Then they are scheduled in a FIFO manner. This discipline is the simplest, most economical and commonly implemented queueing discipline. If the queue length exceeds the available buffer space, the incoming cells are discarded. To minimize cell loss, congestion must be detected by use of two queue thresholds. When the queue length is above the high-threshold level, a congestion indication flag is set initiating marking of packets that will eventually drop the incoming packet rate. It remains set until the queue length drops below the low-threshold level. Single queueing with two thresholds is illustrated in Figure 4-7. Single FIFO queueing, because of its simple nature is easily implemented, particularly at very high-speed ports. It does not, however, provide any local mechanisms to enforce a fair access to buffers and bandwidth, and it leaves such resources open to abuses by malicious flows. In that case, the marking mechanism described suffers in terms of fairness. The problem can be overcome by using processing power to calculate the fair share of bandwidth for each flow and communicate it to its source as in Explicit Rate technique of ABR traffic of ATM. The same external to queueing mechanism must also impose a policing function to punish the malicious sources.

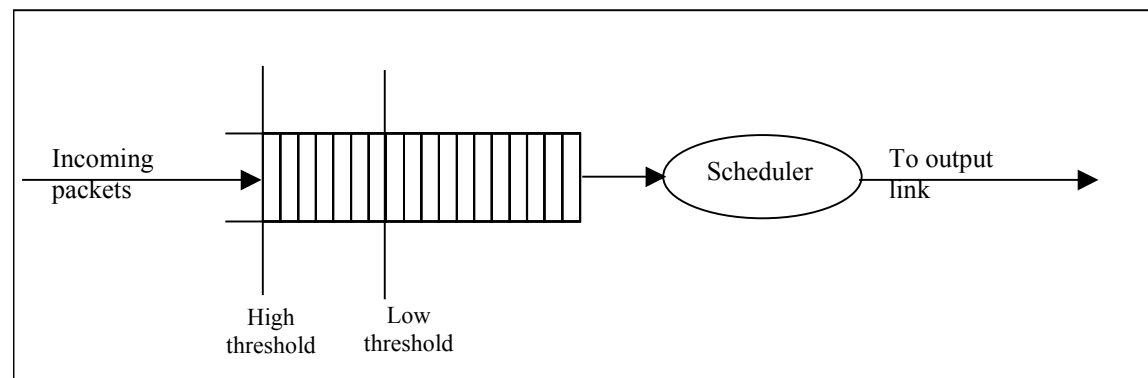


Figure 4-7: Single FIFO queueing and two threshold congestion detection approach

Unlike the first approach, where all flows are queued in one single queue, in the per-Flow approach, cells from different flows (or VCs in ATM networks) are queued in

separate queues and buffer space is allocated in a per-flow basis. Multiple classes of traffic, with different priorities and QoS requirements can be fairly served, with the per-Flow implementation in conjunction with a proper scheduling policy. The output –port buffer space could be divided among all flows in a fixed manner or dynamically shared among them. In the fixed buffer allocation method (static), each VC is only allowed to occupy its own VC buffer share, but in dynamic buffer management, VCs can take more than their share. The state of congestion is determined similarly using the two-threshold approach described earlier. Still two thresholds can be maintained for each of the flows, as seen in Figure 4-8, so that ill-behaved flows can be identified. The isolation provided by the separate queues ensures fair access to buffer space and bandwidth. This also allows the delay and loss behavior of individual flows to be isolated from each other. This per-flow information can be used for congestion control, either by marking of packets (EFCI marking of ATM cells for example), or Explicit rate mechanisms. It can also help identify and police misbehaving sources effectively. In the static buffer management variation, where flows are given a fixed buffer share, the policing can be completely eliminated. This happens because in the static buffer scheme, if a flow is misbehaving, only its queue will grow and overflow. The dynamic buffer scheme, on the other hand, must utilize some intelligent mechanism to achieve policing.

Generally per-VC queuing offers some advantages over single FIFO queuing. These become critical when QoS must be implemented, since the latter when used broadly achieves the most efficient usage of network capacity. Still, per-flow implementation suffers considerably in terms of implementation and scheduling complexity. Since per VC queuing is proportional to the number of maximum VCs that can be served, it doesn't scale well for million of VCs. In addition complex scheduling policies must be implemented on a per-flow basis. When the number of VCs is relatively low (like 64K flows in our IP) buffering complexity is low but the scheduling complexity is still high [1, chapter 9], [8], [9], [10,] [11].

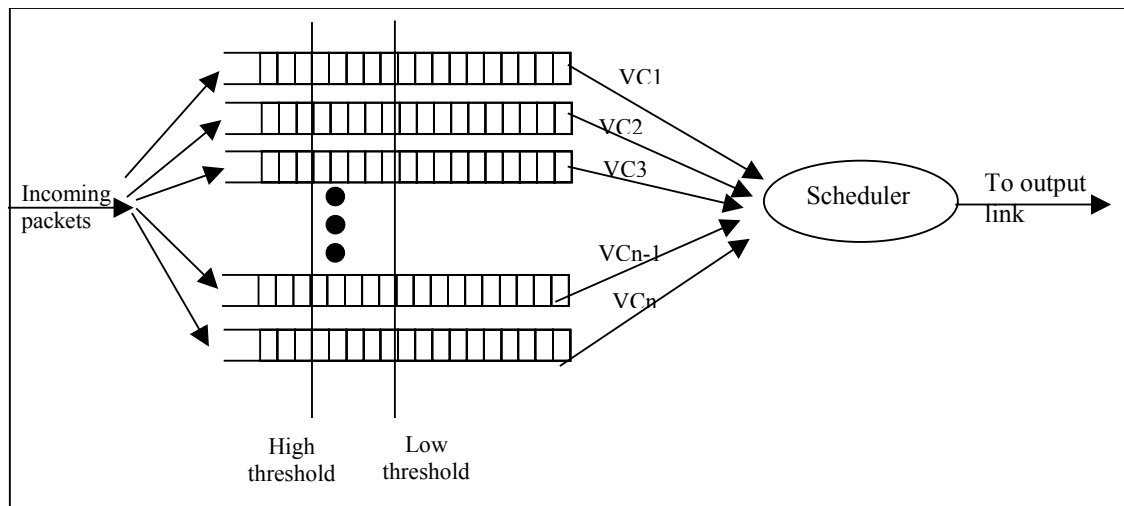


Figure 4-8: Per-Flow queuing and two-threshold detection approach

5 The DIPOLO ATM Switch

5.1 The DIPOLO Architecture

The general architecture of the DIPOLO ATM Switch is depicted in figure 5-1. Based on a Second Generation architecture, this system is made of a number of VDSL Line Cards that connect the End Users with the Switch, one or more ATM 155 cards for the interconnection of the Switch with the central ATM Network through OC-3/STM-1 lines, a CPU Card, an ABR Server Card and a shared Cell Bus.

The exact functionality of each card is given below:

- **Shared Cell Bus:** At system level, the exchange of data and management and signaling information is done by ATM connections. The transmission of such data is done through the use of a shared Cell Bus (CellBus). For redundancy the system uses two such busses. One main and one redundant that it is used in case of failure by the main one
- **ATM 155 Card:** It executes all the functionality of the physical level of ATM, like physical level ending, ATM routing, signaling Q.2931, OAM function etc. The connection of many cards over a backplane incorporates the switch.
- **ABR Server Card:** It executes all the necessary functionality for the acceptance, storage and forwarding of cells that carry ABR traffic. The ATM 155 cards can not handle such traffic because it requires extensive buffering. The core of the functionality is the maintenance of a large number of multiple level queues.
- **CPU Card:** It provides top system maintenance and performs functions like Call Admission Control (CAC).
- **VDSL Line Card:** It implements the interconnection of end users with VDSL modems through wire cables.

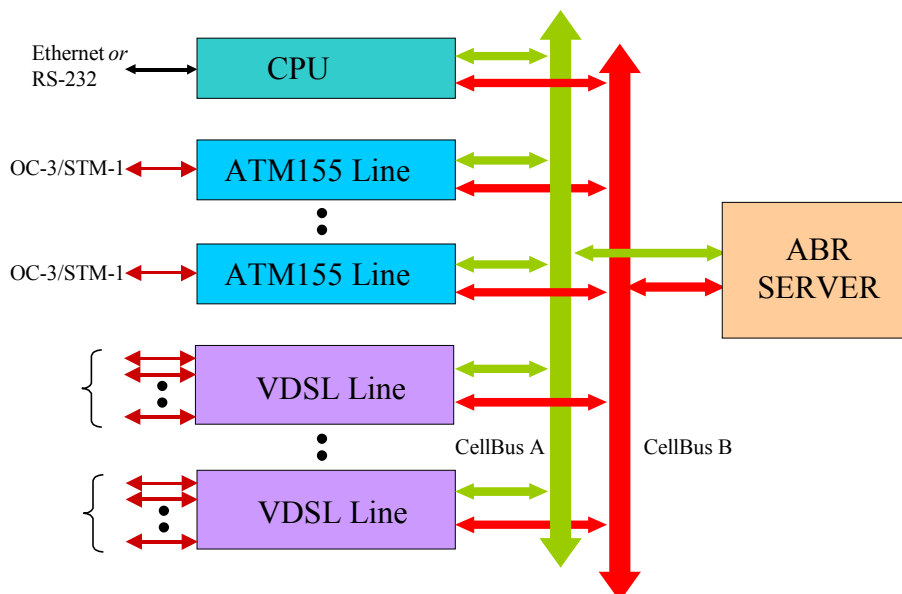


Figure 5-1: General Architecture of DIPOLO ATM Switch

With this organization, the centralized approach of handling the ABR Traffic (by the centralized ABR Server Card), allows the concentration of temporary memory in one card, increases memory utilization and drops the system cost. It also allows the creation of systems with or without ABR support and increases the flexibility of the system in terms of needs and development by the contractors.

5.2 The ATM 155Mbps Card

The ATM 155 Card's architecture is composed by 4 subsystems shown in figure 5-2

- The switching device
- The Cell processing device
- The local control processor
- The local physical level device

In the next subsections we give a description of these subsystems and the commercial products that were used for their implementation.

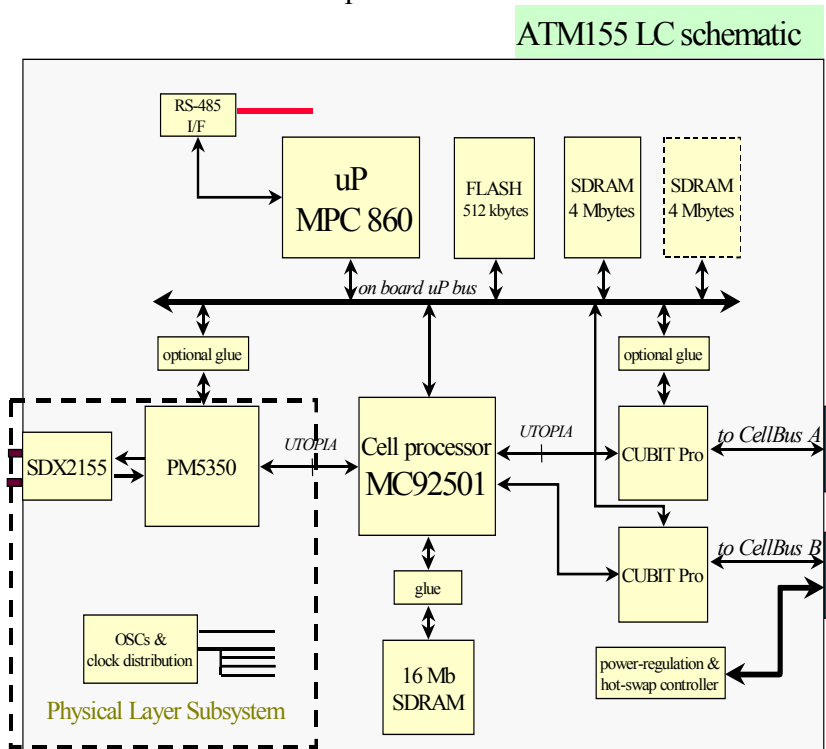


Figure 5-2: The ATM 155 Card block diagram

5.2.1 The switching device (Transwitch Cubit Pro) [18]

The switching device is connected with the cell-processing device through a UTOPIA interface and it comprises the interconnection of the card with the Cell Bus. The latter is located on a separate backplane box that holds together all the cards.

5.2.2 The Cell Processing device (Motorola MC92501 Cell Processor) [17]

For the implementation of the cell-processing device we use the MC92501 Cell Processor by Motorola, that is a well-known and standard processor. It is connected with the physical device and the CubitPro device through a UTOPIA interface. It also connects with the local processor by use of a special processor interface.

The Cell Processor executes all the OAM functions in hardware, and inserts in the system cells that carry signaling and managing information through the local CPU interface. It also executes address translation (VPI/VCI) translation and inserts in the incoming cells routing headers that notify the CubitPro with the card of destination of the cell. The address translation table is located in an external dedicated SRAM that is updated by the local CPU. The MC92501 also performs per flow policing (UPC/NPC) by using a leaky bucket algorithm.

5.2.3 The local CPU (MPC860 by Motorola) [17]

For the local control of the ATM 155 Card, the MPC860 processor by Motorola was used. This processor is responsible for the initialization and normal function of the card. It interfaces with all the devices on the card. The function of the CPU is supported by a DRAM that is used for storing data and an EPROM for storing code.

By interfacing with the other devices the MPC860 can extract information relative to their functionality. It also polls them to get a view of the ATM traffic that passes through the card and collects statistic for system management. It performs signaling while detecting and correcting faults. Special ATM cells that are transmitted over the CellBus through dedicated VP/VC connections do the communication between this CPU and the central CPU card of the switch. Finally the MPC860 uses an RS485 interface for external communication and supervision.

5.2.4 The Physical level device (PM5350 S/UNI-155-ULTRA) [20]

For the STM-1 physical connection the PM5350 S/UNI-155-ULTRA device by PMC Sierra was used. This is a device that interfaces with the MC92501 through a UTOPIA interface and implements the ATM transmission convergence for the SONET/SDH 155.52 Mbit/sec. The PM5350 is used for the implementation of basic functions of the physical level for an ATM UNI (User Network Interface) interface. In a typical STM-1 application, the device performs clock and data recovery at the reception side and clock generation at the transmission side. For the optical transmission of data an SDX2155 simple rate optical transceiver is used for SDH STM-1. It can support a connection length longer than 15 km and is manufactured with 155 Mbit/s standards set by the ATM Forum. Finally the SDX2155 has as PECL interface for communication with the PM5350 device.

5.3 The CPU Card

The basic functions performed by the CPU Card are:

- Central administration and monitoring of the DIPOLO Switch
- Call Admission Control (CAC)
- Interface with external management systems (support for SNMP and/or CMIP)

In figure 5-3, the block diagram of the CPU Card is given.

The basic parts of the Card are described in the following subsections.

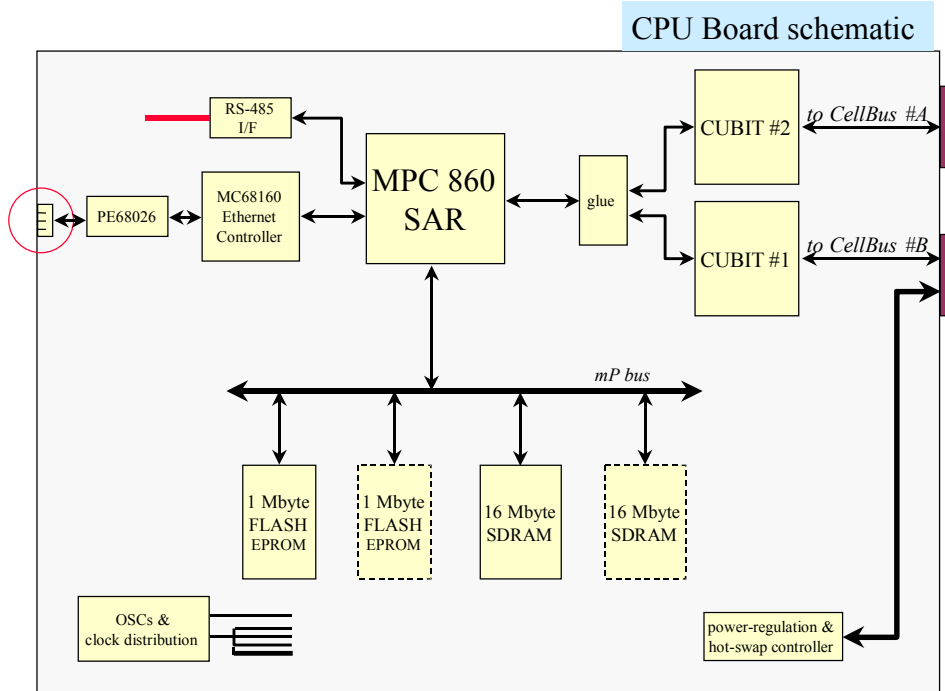


Figure 5-3: The CPU Card block diagram

5.3.1 Motorola MPC869SAR (PowerQUICC)

The MPC860SAR is advanced version of the Motorola's MPC860 processor, based on the PowerPC architecture. It executes additional ATM and SAR (Segmentation and Re-assembly) functions. One of its basic advances, in terms of system design, is that it provides interfaces to memories, serial transceivers and bus buffers. The MPC860SAR provides the following interfaces:

- Interface with SRAM, DRAM, EPROM, Flash and other peripherals.
- Seven serial interfaces controlled by an integrated communication subsystem processor that can support a number of different communication protocols.
- Half-duplex UTOPIA interface at 155 Mbps that can support multi-PHY function and extended ATM cell size.

The processor performance reaches 52 MIPS at 40 MHz, with cell processing (incoming and outgoing) rate of over 60Mbps.

5.3.2 The CubitPro Device

Like on the ATM 155 Card, the Cubit Pros provide the interface of the CPU Card with the double Cellbus, allowing the communication of the CPU (MPC860SAR) of the CPU Card with the local CPUs of the other cards, with the use of Control Cells.

5.3.3 Memories

Memories used are:

- SDRAM: 16Mbytes for the storage of code and data
- ROM: 2MB
- Flash EPROM for the storage of code for the Central and the local CPUs.

5.4 The ATM Line Card

The E1 Physical Layer part terminates 16 E1 cell-based lines. All terminated data are interchanged with the MC92501 cell-processor over a UTOPIA II bus. As shown in figure 5-4, the E1 physical layer sub-system consists of the following components:

- 4 PM7344 QUAD E1/T1 Framers
- 4 PM4314 QUAD E1/T1 Line interface components

In figure 5-4, the block diagram of the VDSL Line Card is given.

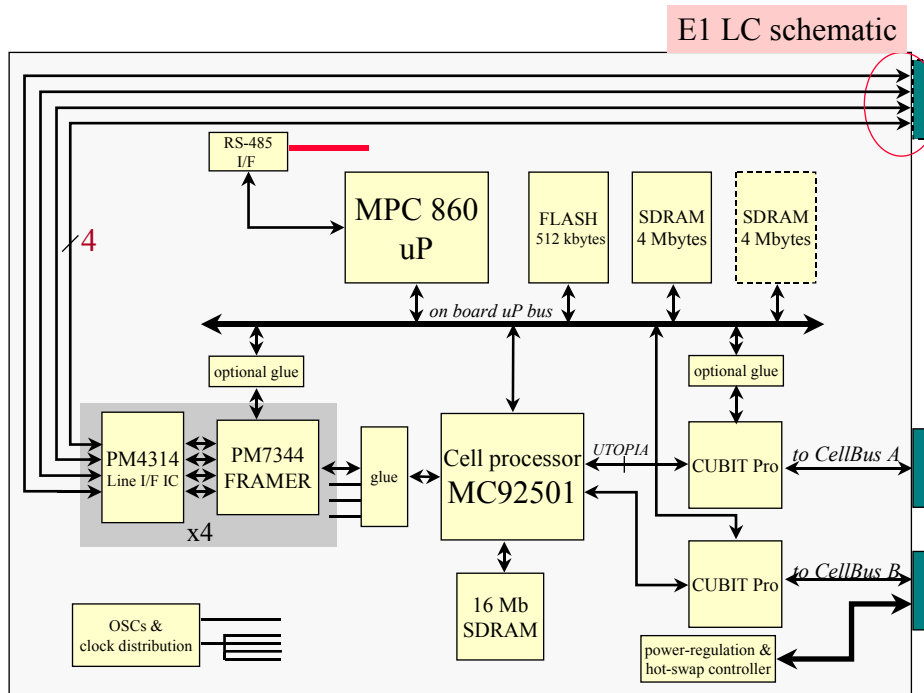


Figure 5-4: The VDSL Line Card Block Diagram

The incorporated parts are described below.

5.4.1 Framer part - PMC-Sierra PM7344 [20]

The PMC-Sierra PM7344 is a 4 port (input and output) Framer device. It accepts raw line data (RZ encoded) from the Line interface circuit (PMC-Sierra PM4314), and provides cells to the ATM Board through a UTOPIA 2 like interfacing (SCI-PHY). Some glue logic (FPGA) is required to convert the SCI-PHY signals to pure UTOPIA 2 signals. The PM7344 provides cell scrambling (from UTOPIA to PM4314), de-scrambling (from PM4314 to UTOPIA), and cell delineation (according to G.804). It inserts/extracts idle-unassigned cells to/from the data stream. Rate de-coupling of the UTOPIA interface is achieved by a 4-cell input FIFO and a 4-cell output FIFO. The PM7344 is connected to the Motorola MPC860 with an 8-bit CPU interface and the appropriate glue logic. The MPC860 assign values to the internal registers of the PM7344 in order to program it, and gets information about the line status by reading the appropriate registers. The PM7344 can also generate interrupts to the MPC860 under certain line conditions.

5.4.2 Line interface circuit part - PMC-Sierra PM4314 [20]

The PMC-Sierra PM4314 integrates 4 duplex E1 (G.703) compatible line interface circuits. It accepts E1 (G.703) bipolar line signal from the electrical components, and provides digital (RZ encoded) pulses to the Framer part (PM7344). The PM4314 provides clock recovery and performance monitoring in the receiver. It is also equipped with a generic microprocessor interface (connected to the MPC860 with the appropriate glue logic) for initial configuration, ongoing control and status monitoring. The microprocessor interface utilizes an 8-bit data, and a separate 8-bit address bus, and is able to generate interrupts upon detection of various alarms, events, or changes in status.

5.5 The ABR Server Card

The ABR ATM Traffic covers the needs of traditional LAN networks. Computational systems on a LAN want to send data, the minute these are generated, at the highest possible rate, but without congestion that causes cell loss. The reason behind this is that computer data are sensitive to loss (as opposed to multimedia data like voice and video where loss can be tolerated to a point), and retransmissions can seriously degrade the network performance.

Instead of reserving network resources for the bursty traffic of LANs, it can be served by the ABR service that uses the available bandwidth, left unused by other sensitive to delay and jitter, traffic services of ATM networks, like CBR and VBR. In order though to allow for end-users to send whatever they want whenever they want without loss, switches with extensive buffering capabilities should do the service of ABR traffic. In that way, cells that can not be transmitted during time of congestion can be stored for some time.

For this reasoning, the DIPOLLO Switch utilizes a centralized ABR Server card, in order to efficiently utilize the memory resources available at the lowest of costs. The

ABR Server Card can provide for the buffering of all the ABR traffic of the Switch. In that way, supplying memories to all the card thus increasing complexity and cost, is avoided. Memory utilization is also the highest. The downside is that some CellBus bandwidth is lost since each ABR traffic cell traverses the bus twice (once from card of entry to ABR Server card and once from ABR Server card to card of exit) and that one ABR Server card is a single point of entry. Since though the switch is designed for low-end home application, where multimedia data are the core rather than computer data, and cost matters for a the highest possible market penetration, centralized buffering of ABR Traffic is justified.

The architecture of the ABR Server Card is given in figure 5-5. The card is made out of the following devices:

- The Switching device (Transwitch Cubit Pro)
- The ABR Server Unit (ABRSU) (EPF10K200EBC600-1 by Altera)
- The Memory module (256 MB of SDRAM)
- The local processor (MPC860 by Motorola)

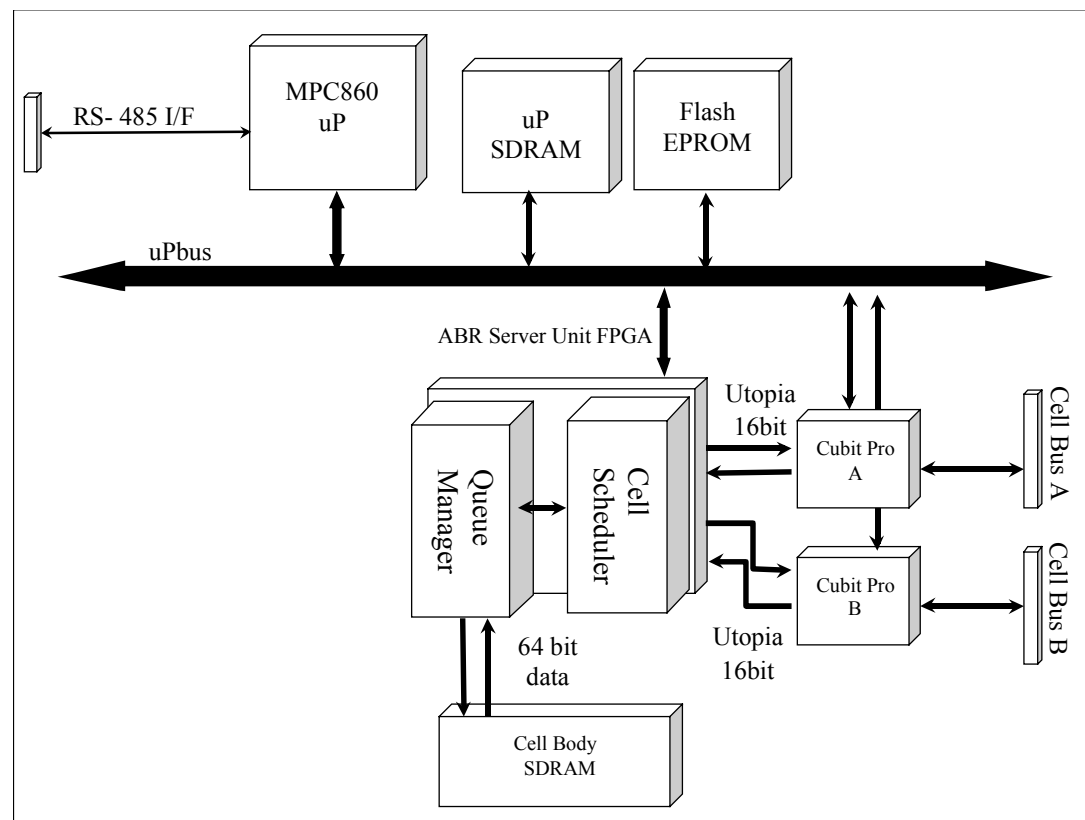


Figure 5-5: ABR Server Card block diagram

5.5.1 The Switching device (Transwitch Cubit Pro)

See Section 5.2.1.

5.5.2 The ABR Server Unit (ABRSU) (EPF10K200EBC600-1 FPGA by Altera) [16]

The ABR Server unit (ABRSU) is composed by two main blocks, one for buffering incoming cells, the Queue Manager and one for scheduling their departure, the Cell Scheduler. The ABRSU interconnects with the CubitPro Switching devices through a two way UTOPIA interface 16 bit wide. It also has a CPU Interface to communicate with the MPC860 local processor. Finally it has an SDRAM controller to access the buffer space of the SDRAM memory module that stores the Cell Queues and the queuing data (head tail pointers for the queues and other).

The purpose of the ABRSU is to store ABR cells according to the 16bit field that is in the header of every cell and defines uniquely the Virtual flow (VP/VC) that it belongs. The ABRSU must also schedule the transmission of each cell. The whole logic of the ABRSU is implemented by an EPF10K200EBC600-1 FPGA by Altera. The ABRSU can support up to 64K flows of ABR traffic, since the routing ID is 16 bit wide.

5.5.3 The Memory module (256 MB of SDRAM) [19]

A 256MB DIMM of SDRAM was used as the dedicated buffer of the ABRSU, where cells were stored in queues as well as other queue data (head and tail pointers for the queues) and per flow information. Although less than 256 bytes could suffice for the application, 256Mbytes of SDRAM are common DIMMs on the market.

5.5.4 The local processor (MPC860 by Motorola)

See Section 5.2.3

6 The ABR Server Unit (ABRSU)

In this chapter we give a presentation of the ABRSU that includes the subject of this thesis, the Queue Manager. It also includes additional blocks that interface with the Queue Manager in order for the latter to implement the Queueing function for the ABR Server card of the DIPOLO Switch. The most important is scheduling block, the Cell Scheduler that is subject of this thesis.

6.1 The ABR Server Architecture

The ABR Server Unit is implemented on the ABR Server Card by Altera's EPF 10K200EBC600-1 FPGA. It provides the following functions to the Card:

- Acceptance of Incoming ABR traffic cells through a dedicated UTOPIA input interface with the CubitPro devices of the card.
- Recognition of RM cells so that RM marking can be performed on them in case they do not conform to their quality of service.
- Storage of ABR,RM cells in the dedicated cell memory (SDRAM DIMM) based on the flow they belong. (per flow queueing). Tail pointers of the queues are also enumerated with the new arrivals.
- Grouping of the flows into group of flows (Flow Groups). The grouping is free of any constraint and can be used to group flows either by the Dipolo Switch port/card of destination or the Quality of Service they have assigned with or both.
- Scheduling of the transmissions of stored cells that lay on the head of the flow queues from the ABR Server Card. The transmission is based on the available bandwidth of the CellBus. For this function, a Scheduling block has been implemented that uses the information of the quality of service parameters for each flow group to request the dequeing of one cell from a group of flows to be transmitted. That cell will then be dequeued from the SDRAM by the Queue Manager block that will enumerate the head pointer of that flows queue and send the cell for transmission. When the correct transmission to the card of destination is acknowledged then the memory block that the cell was kept will enter a free list and will eventually be used by another incoming cell [23].
- Transmission of the dequeued cells to the CubitPro devices of the card through a dedicated UTOPIA output interface. There is also the capability of retransmission of the cell in case the first transmission over the CellBus by the Cubit fails.
- Provision of a CPU Interface block that is used for interconnection of the ABRSU logic with the MPC860 local processor. Through this interface the MPC860 can initialize or terminate the ABR flows that are served by the card. It can also have access to the data structures of queues and the queues themselves that are stored in the SDRAM. The MPC 860 can also access the Flow Control parameters of each Queue dynamically by providing the thresholds that are used for marking of flows individually. It also can access the service times used by the Scheduling device to request the service of Flow Groups.

In figure 6-1 the block diagram of the ABRSU is given:

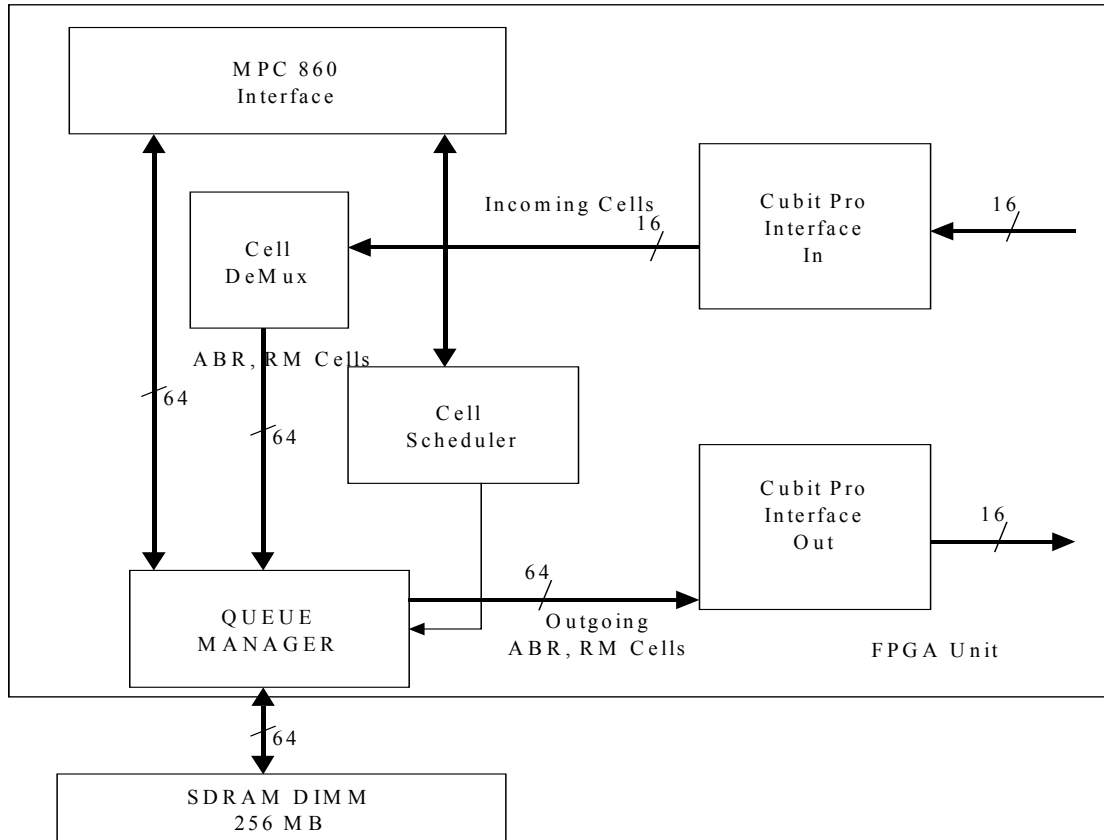


Figure 6-1: The ABRSU internal block diagram

The internal sub-blocks of the ABRSU are presented in the following chapters, while a separate chapter is dedicated for the Queue Manager IP that is the subject of this thesis (See Chapter 7).

6.2 The CPU Interface

The CPU Interface subblock, as mentioned previously, is responsible for the communication, initialization, and dynamic configuration of the rest subblocks of the ABRSU. Through this subblock the CPU initializes the Flow Data Structures stored in SDRAM by issuing special commands executed by the Queue Manager subblock. If needed the Queue Manager will return data to the CPU through this sub-block. The CPU can also dynamically configure QoS service parameters at the Cell Scheduler [22],[23].

The internal organization of MPC 860 Interface Block is given at Figure 6-2. Since the MPC860 bus on the ABR Server Card works at clock speeds ranging from 10 to 25 MHz while the FPGA can achieve a higher rate internal clock, the MPC860 Interface subblock must resign on both these clock domains and use some kind of synchronization.

From the CPU side the MPC Interface Subblock uses the CPU bus pins as input/output. These bits is a 32 bit wide data bus, plus some control signals (chip select, Read/Write etc). Through these pins the CPU writes at special configuration registers and configures the Queue Manager and the Cell Scheduler. Μέσω αυτών των pins ο Μικροελεγχτής μπορεί να γράφει και σε συγκεκριμένους καταχωρητές και

έτσι να διαμορφώνει δυναμικά την λειτουργία και τις δομές των υπομονάδων Queue Manager Block και Cell Scheduler Block. It can also receive data by reading dicated for this purpose registers.

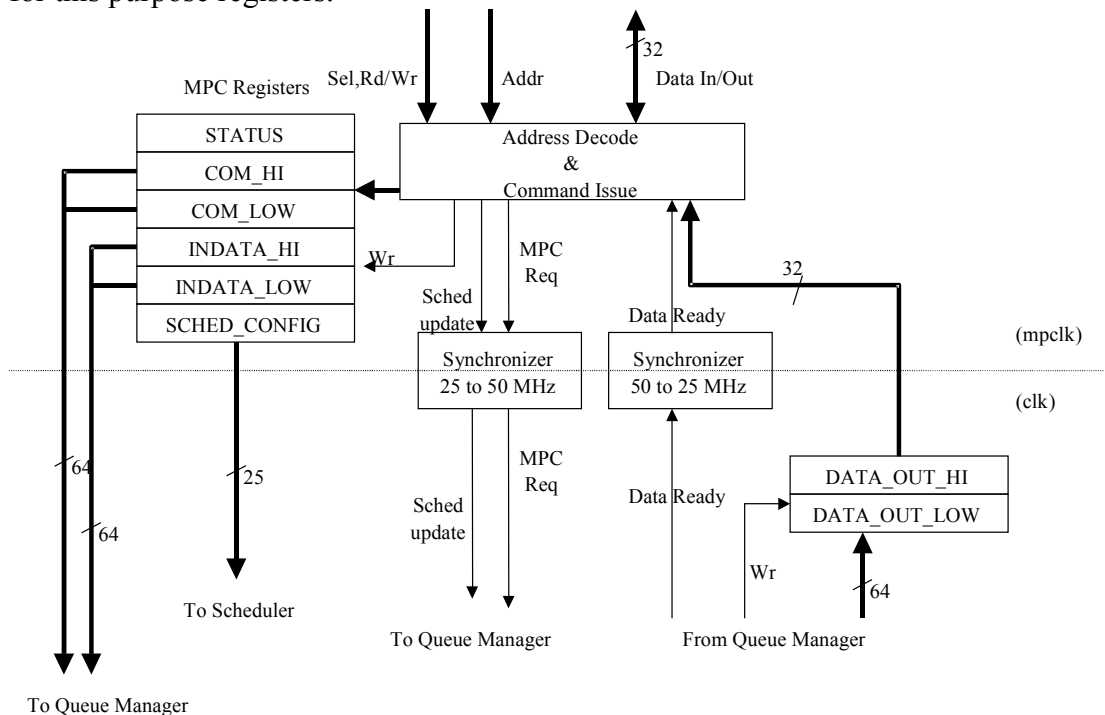


Figure 6-2: The CPU Interface sub-block diagram

The subblock seen in figure 6-2 that is responsible for the access of the CPU dedicated registers is Address Decode block and Command Issue. Its purpose is to decode the data of the CPU address bus and execute a Read or Write Operation according to the address. When special registers used for command request are written by the CPU, a signal named Command Request is generated by the subblock to notify (after synchronization) the Queue Manager that a command has been issued by the CPU, that must be executed by the Queue Manager. The same function has the signal Sched Update that is generated each time the CPU writes the dedicated QoS register. This signal after synchronization notifies the Cell Scheduler of the change.

6.2.1 MPC 860 Interface Block - Queue Manager Block Interface

The MPC 860 Interface Block can configure the Cell and Queue data Memory unit that lies in the SDRAM. The SDRAM fully controlled by the Queue Manager and CPU can access it by placing certain commands to the latter. Some command include data to enumerate the Queue and flow structures, others aren't, while others return data to the CPU.

The commands that the CPU can place at the MPC860 Interface block concerning the Queue Manager and the registers used for each (Refer to figure 6-2) are the following:

- 7) **Open Flow:** Initializes flow, requires write access of registers COM_HI, COM_LOW.
- 8) **Close Flow:** Terminates a flow, , requires write access of registers COM_HI, COM_LOW.

- 9) **Write:** Writes a 64-bit word to SDRAM. It requires write access of registers COM_HI, COM_LOW, INDATA_HI and INDATA_LOW.
- 10) **Read:** Reads a 64-bit word to SDRAM. It requires write access of registers COM_HI, COM_LOW.
- 11) **Read Counter:** Reads the counter of cells queued for a certain flow. It requires write access of registers COM_HI, COM_LOW
- 12) **Change Parameters:** Change the High and low thresholds used for cell marking of misbehaving flows. It requires write access of registers COM_HI, COM_LOW.

For each of the commands to be issued, at least two registers must be written. It is also necessary, that the last register to be written should be COM_LOW, since it sets signal Mpc_Req. The latter after passing synchronization from 25MHz of MPC 860 Interface Block to 50 MHz of FPGA logic, notifies the Queue Manager sub-block, that there is a new Microprocessor command to be executed.

Read and Read Counter Commands request data to be returned to the CPU by the Queue Manager IP. When the latter creates these data, stores them in registers DATA_OUT_HI and DATA_OUT_LOW. The two 32 bit are necessary since the data are 64 bits long and the MPC860 Interface is only 32 bits long. DATA_OUT_HI stores the 32 most significant bits while DATA_OUT_LOW the least significant. Writing of these two registers sets the signal Data_Ready that after passing synchronization from 50 MHz to 25 MHz it set the corresponding bit in the STATUS register. The CPU polls the status register and learns that the data are ready for accessing. The STATUS register also informs the CPU that it is able to issue a new command by another bit.

The command arguments and their execution by the Queue Manager is presented in detail in subsection 7.2.4.

6.2.2 MPC 860 Interface Block - Cell Scheduler Block interconnection

The MPC860 Interface Sub-block allows the MPC860 to shape the Quality of Service given to each one of the 64 Flow Groups that can exist at any time at the ABRSU, by programming the Cell Scheduler. Writing QoS parameters to the SCHED_CONFIG register does this. In this register the CPU writes the minimum service interval which is the minimum time between two successive dequeues from flows in that Flow Group- and the id of the respective Flow Group. When this register is written by the CPU, the pulse signal, Sched_update, is set and after passing from the synchronization logic it is used as a write enable for the Service Interval Memory that exist in the Cell Scheduler.

6.3 The Cell Demultiplexor

The Cell De-multiplexor block is responsible for the accumulation of a hole cell in 7 words of 64 bits each, in order to be sent to the Queue Manager sub-block for enqueueing in the appropriate queue. This block is necessary because first enters in the Cubit Interface block and stored in a 16-bit fifo and must be restored in 64-bit fifo so that when enqueue operation is performed successive 64-bit word should be provided to the SDRAM Interface on each clock cycle.

The Cell Demux FSM polls the avail signal coming from the Cubit Interface sub-block to see if there is a cell for enqueue. If this signal is one then a whole cell is waiting in the Cubit Interface blocks FIFO. The FSM then proceeds in filling a 64-bit register with 16 words read with the appropriate signal from the Cubit Interface block. Each time the register fills a write operation is done by the FSM to the Input cell buffer with the data of the 64-bit register. The Input Cell Buffer is a 7x64 show ahead FIFO (show ahead FIFO is the one that outputs immediately the first datum after a write). When this FIFO is full a cell is ready to be enqueued and actually the `fifo_full` signal of this FIFO is used as an enqueue request to the Queue Manager. No new cell import can begin by the FSM until the FIFO becomes empty again. This is done by the successive reads by the Queue Manager from this FIFO. The first datum of each cell contains the Flow ID that the Queue Manager uses for enqueueing the cell to the respective queue (This is the reason that the FIFO is show-ahead). After the Input cell buffer empties, a new cell can be loaded into it and wait for its enqueue operation. The Cell Demultiplexor sub-block can be seen in Figure 6-3.

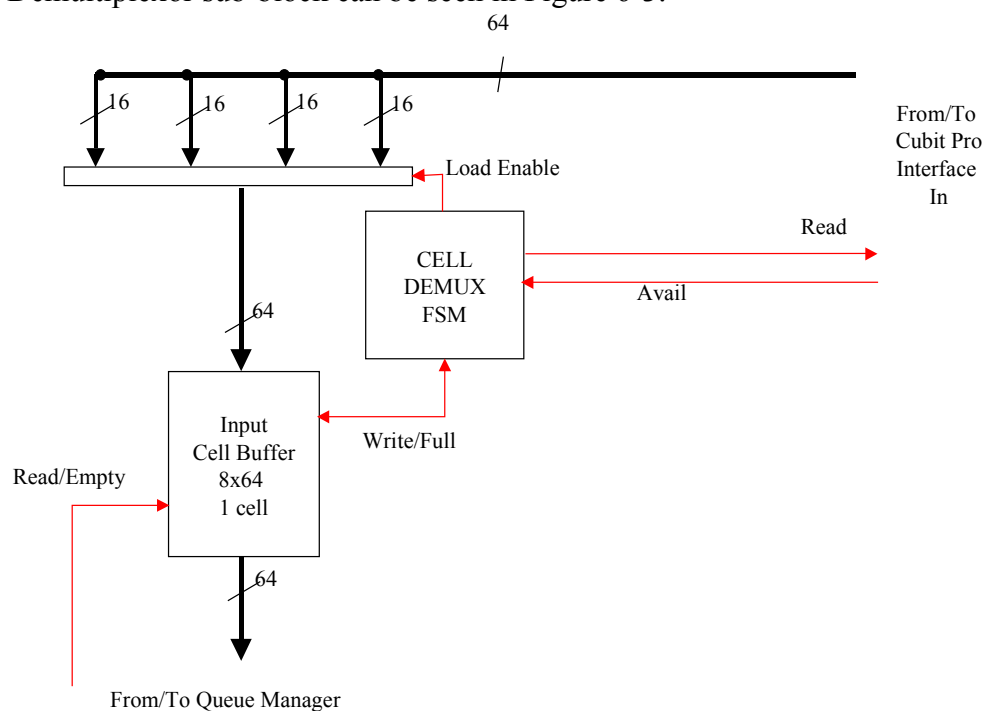


Figure 6-3: The Cell Demultiplexor sub-block diagram

6.4 The UTOPIA Interfaces

There are two UTOPIA interfaces, both of 16 bits, incorporated in the ABRSU. The first is called UTOPIA Input Interface and is used to import the cells from the Cubits of the ABR Server Card into the ABRSU in order to be enqueued. The second is called UTOPIA Output Interface and is used to sent the dequeued cells from the ABRSU to the Cubit devices of the ABR Server card for transmission over the CellBus to the card of destination. A brief description of these two interfaces is given in the following two subsections.

6.4.1 UTOPIA Input Interface

A simple block diagram of this 16-bit interface is given in figure 6-4. The Cubit Devices send data to the FPGA (ABRSU) through a 16 bit Utopia interface. This block has a controller that recognizes and responds to Utopia signals. It can therefore be notified of the beginning of a cell and store its 16-bit words in an elastic FIFO, with the use of the utopia clock provided by the Cubit. This FIFO is 8 cells long. A separate counter exists to count the numbers of cells that exist in the FIFO. Note that a cell is considered to be inside the FIFO (Counter ++) only when the whole cell has entered the FIFO. This is done, in order to allow the Cubits to send a cell inconsecutively and the rest of the ABRSU to receive it consecutively. If the FIFO is full then the controller notifies with the avail signal to send no other cell. On the other side the Cell Demultiplexor see an avail signal that notifies it (After passing synchronization from the Cubit Clock to the FPGA Internal Clock) that there is a cell inside the FIFO. The Cell Demultiplexor then proceeds in reading the cell out of the FIFO and into the Input Buffer FIFO, by use of the read port of the FIFO that has different clocking than the write port. The port clock is FPGA internal clock. Cell Counter is decreased when the whole cell has been read out of the FIFO.

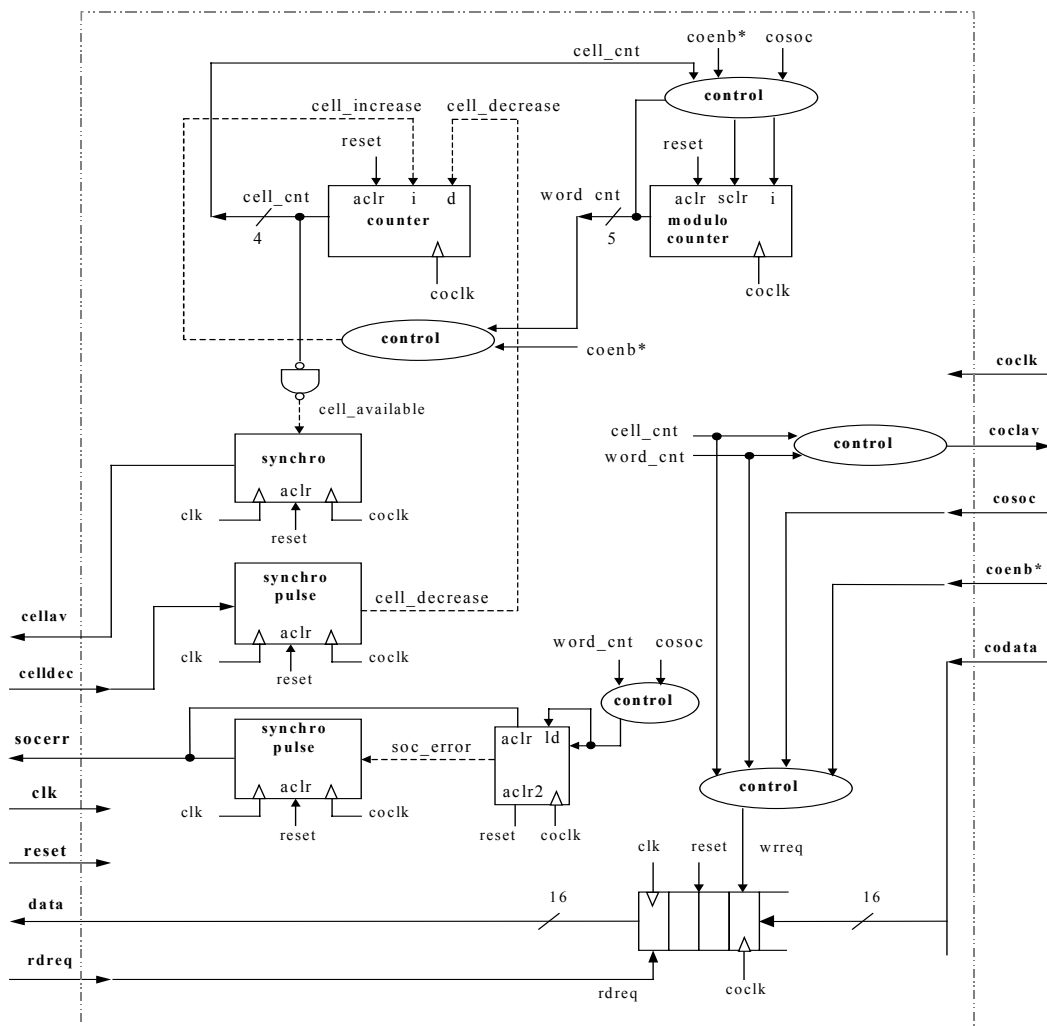


Figure 6-4: UTOPIA Input Interface block diagram

6.4.2 UTOPIA Output Interface

The UTOPIA Output Interface is responsible for forwarding the dequeued cells from the ABRSU to the Cubits for transmission over the Cellbus. From the side of the Cubits it operates on the UTOPIA clock, produced by the Cubit Devices. A Controller sub-block uses this clock along with the UTOPIA handshake signals to transmit 16-bit word of a cell stored inside an elastic FIFO. This FIFO 8 cells long. If the Cubit internal FIFOs are full then the controller stops transmitting new cells. On the Side of the ABRSU, the Cell Scheduler that receives a cell after dequeuing uses the write port of the FIFO that operates on the FPGA clock to fill it. If the cubit accept no cells then the FIFO becomes full and the Cell Scheduler receives an no apace available signal. The scheduler then stops scheduling any new transmissions. An external to the FIFO counter is used to count the number of cells inside the FIFO. Signals that increase this counter, as well as avail signals, pass from synchronizing logic that lies between the clock domains.

A block diagram of the UTOPIA Output Interface is given in figure 6-5.

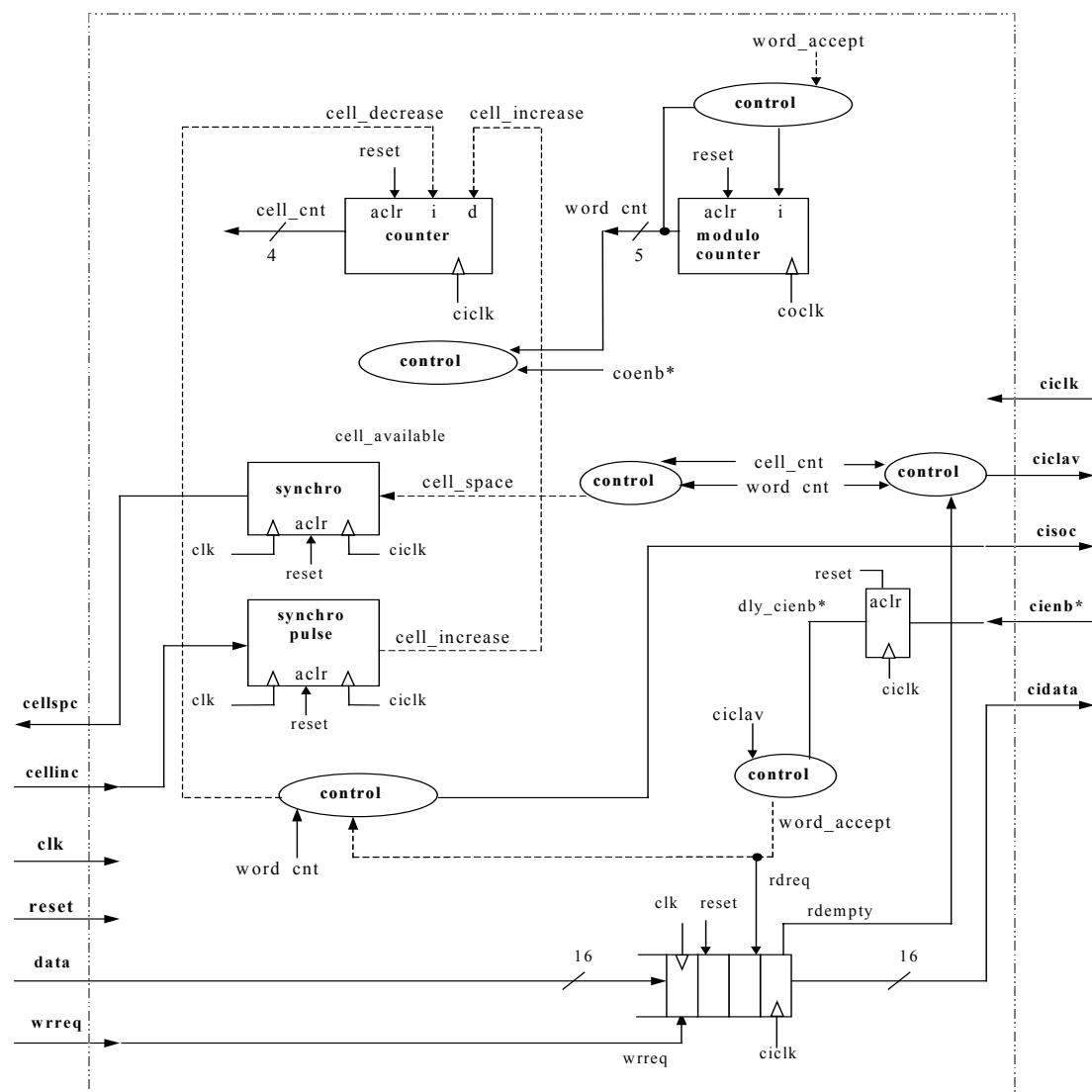


Figure 6-5: UTOPIA Output Interface sub-block diagram

6.4.3 The Pulse Synchronizer

Figure 6-6 gives the internal organization of Synchronizing block that is used to pass signals between the UTOPIA interface clock domain and the rest of the ABRSU clock domain.

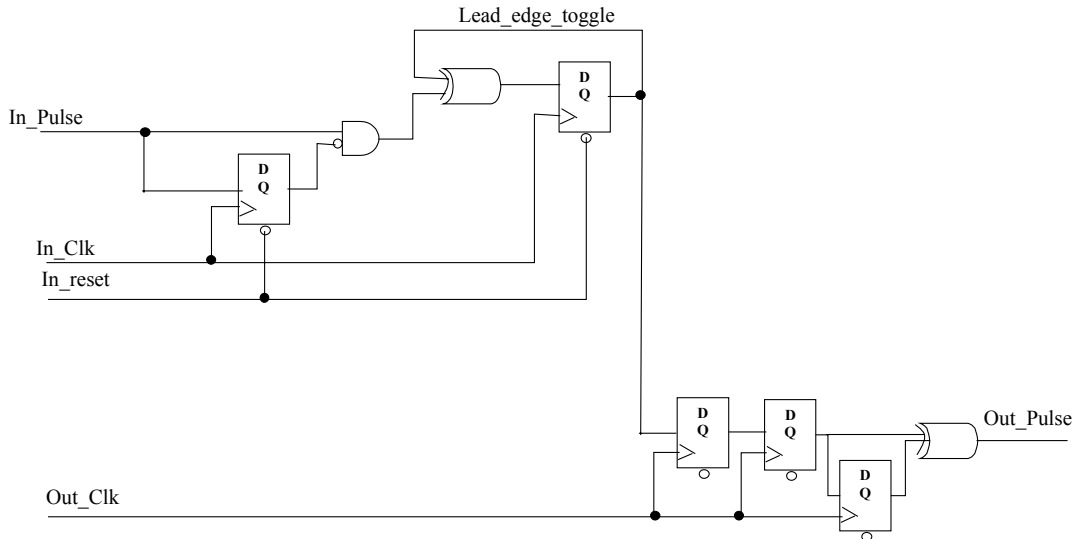


Figure 6-6: The Pulse Synchronizer

The main features of this synchronizer are:

- **It requires the following to be true:**
Cycles between in pulses $> 2 * T_{out} / T_{in}$
- **it guaranties that there is more than two output clocks between changes of lead_edge_toggle**
- **no conditions on clock relation, frequencies duty cycles and pulse widths**

6.5 The Queue Manager

This block is the main subject of this work and is described in detail in Chapter 7.

6.6 The Cell Scheduler

The Cell Scheduler is, along with the Queue Manager, one the two most important sub-blocks of the ABRSU. While the Queue Manager implements the queuing architecture of the ABRSU and the ABR Server Card in general, the Cell Scheduler implements the scheduling policy. This block can distinguish 64 separate flow groups each with different QoS characteristics and services them by requesting by the Queue Manager to dequeue a cell from a flow belonging to the group for transmission over the CellBus. Flows in each of the Flow groups are serviced equally with a Round Robin mechanism. Maintaining the structure of each Flow Group is a responsibility of the Queue Manager. The main functionality of the block is:

- Maintain a memory with Service Intervals for each of the 64 Flow Groups. These service intervals give the minimum time between subsequent dequeue operations on a Flow Group. The CPU, through the MPC860 Interface register SCHED_CONFIG writes these service intervals to the Service Interval Memory.
- Schedule the next request for a dequeue operation from a specific Flow Group and send it to the Queue Manager
- Receive the cell from the Queue Manager along with the address of the block inside the SDRAM where the cell is stored.
- Forward the cell to the CUBIT Output Interface for transmission over the Cellbus. Wait for an acknowledgement from the Cubits that the cell was transmitted successfully.
- If a negative acknowledgement arrives (the cell was transmitted with an error) log the reason in a special memory that notes which Flow Groups have congestion on their destination Cards and perform exponential back off on the retransmission of the failing cell. When the time of retransmission arrives, request from the Queue manager to read again the cell from the SDRAM, by providing the first with the address of the cell inside SDRAM memory, kept from the first dequeue.
- If a positive acknowledgement arrives then schedule the next request from the Queue Manager and place the address of the cell just transmitted to a FIFO of free cells that can be reused. This FIFO is read by the Queue Manager when is not empty, to get a pointer to a free block that can be used by the first to enqueue a new arriving cell. This saves the Queue Manager from an access to its list of free pointers and thus an access to the SDRAM. The whole procedure is called free list bypassing.

The cell scheduler is not discussed further since it is not the subject of this work but more can be found on [23]

7 The Queue Manager IP

In this chapter, the Queue Manager IP is presented. The IP was implemented with the use of Verilog Hardware Description language and synthesized by MaxPlus II synthesis tool, to fit in an Altera FPGA. This IP uses only one SDRAM DIMM as the only memory module to store the incoming ABR cells, as well as metadata for the logical queues of cells, thus dropping the pin count and costs of the design. Per-flow queueing is implemented as the queueing architecture. This means that each separate flow of traffic served by the switch reserves a special FIFO queue for storing its cells. The number of maximum flows of traffic that can be served is 64K at any time that is more than enough for the needs of an edge switch. A list of free buffers is also maintained, to keep all the free buffers that are not used and provide them for storing incoming cells. Thus dynamic memory allocation is implemented allowing specific flows to have queues of arbitrary size and for better memory utilization. Flows of traffic are also organized in Flow Groups by using cyclic lists. Metadata for these flow groups lists is kept in SRAM memories internal to the FPGA. Initialization of flows and definition of the flow group that they belong is being done by the card micro-controller at connection setup, with set up commands issued through the CPU interface of the IP. Flow Control parameters can also be set during flow initialization by the CPU. It does this, by setting information on the maximum queue size allowed to the respective flow. Beyond this size, EFCI and RM marking as described by the ATM forum is performed.

Since only one memory module is used for the needs of the IP, no parallel accesses can be made. Only sequential accesses are possible during the enqueueing or dequeueing of a cell. This fact makes the memory the bottleneck of the queueing bandwidth. In order to reduce memory accesses buffer pre-allocation to each of the 64K queues is implemented, as well as free-list bypassing with the help of logic external to the IP (the Cell Scheduler in the ABRSU environment). The SDRAM controller that is incorporated in the IP to handle the accesses with the SDRAM module is also carefully designed to handle consecutive accesses without loss of clock cycles. The enqueue and dequeue operations are designed in a way that avoids data dependencies (the result of a read operation is used as the address of the next read or write access) that can induce idle clock cycles during the execution.

Except the interface with the CPU and the SDRAM module, the IP keeps an interface for receipt of cell to be enqueued and an interface with the scheduling module (the Cell Scheduler in the DIPOLO environment) that makes requests for dequeueing cells from a specific flow group. The request is served by the IP dequeueing a cell from that flow group and sending it to the interface in 64 bit data words.

7.1 The Queue Manager IP Architecture

In Figure 7-1 a sub-block diagram of the Queue Manager IP is given.

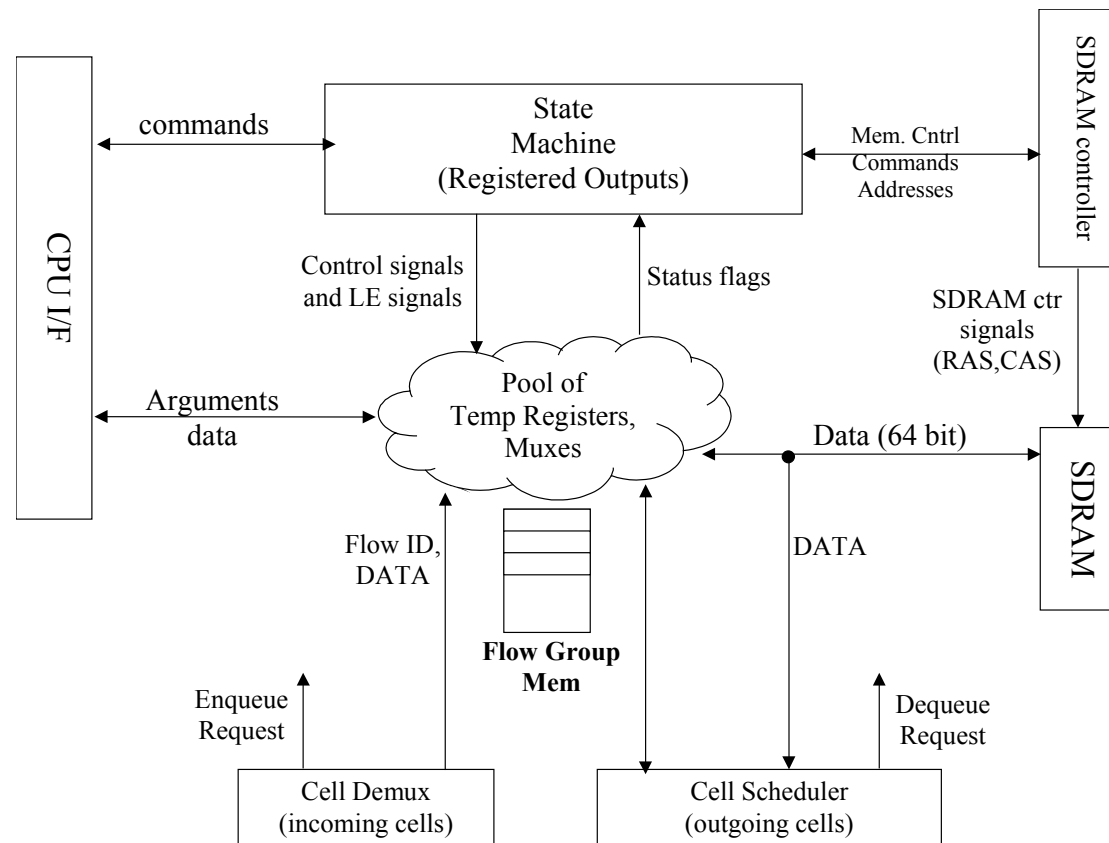


Figure 7-1: The Queue Manager IP sub-block diagram

As seen on the figure the Queue Manager IP is composed by:

- **A Large State Machine:** This state machine is actually composed of many state machines each dedicated to controlling the rest of the logic and memory of the IP in order to implement a specific command. All of these State machines are initiated by a TOP state machine that arbitrates which command should be executed at a given time.
- **A Pool Of Temp Registers and Muxes (Datapath):** These temporary registers are loaded with data coming from all the blocks that interface with the IP. They are loaded with load enable signals produced by the State Machines and they are preceded by muxes that select the origin of the data stored. Again, the State Machines sets the select signals of these muxes. Names of these registers and their respective muxes, to name a few, are: Flow ID, Head Pointer, Tail Pointer, Cell Counter, Hi Watermark, Flow Group ID (of a certain queue accessed at a given command), Free List Head, Free List Tail, Free List Counter (Queue information for the list of Free Buffers of the Cell Memory (SDRAM)) etc.
- **Flow Group Memory:** This Memory keeps data for the 64 Flow Groups supported by the Queue Managers. It is a 64 word memory that stores the Head of

the list of active flows that belong in the Group, the tail of the list, and the status of the Flow Group (if the Flow Group has any active Flows or not). A Flow is considered active when a cell of this flow exists in the Cell Memory (SDRAM) and must be dequeued in the future. The State Machines are doing control of this memory.

- **The SDRAM controller:** This block implements the accesses of the SDRAM DIMM on behalf of the State Machines. It sets the SDRAM DIMM control pins and programs it to read/write bursts of 1,2,8 words of 64 bits each. Data are given directly by the rest of the IP in synchronization with the control signals of the sub-block. No synchronization is implemented since the FPGA that implements the IP and the SDRAM DIMM use the same carefully distributed clock.

The Queue Manager IP interfaces with the following blocks:

- **The CPU Interface:** This interface allows a given CPU the ability to configure the IP, initialize flows, configure the flow control attributes of flows and the QoS parameters of Flow Groups and extract debug information.
- **The Cell Demultiplexor Interface:** By this Interface the IP accesses a 64-bit FIFO that contains a cell to be enqueued. If the FIFO is full a cell exists and the full signal is used as an enqueue request to the State Machines.
- **The Cell Scheduler:** This is the interface with the scheduling block that implements the scheduling architecture. A dequeue request can be given to the IP through this interface along with the ID of the Flow Group that the dequeued cell must come from. The IP responds by sending the cell in 64-bit words, along with the ID of the Flow it belonged to and the address of its respective buffer in the SDRAM DIMM. When free-list bypassing is implemented, such addresses are given back to the Queue Manager for restoring newly arrived cells.
- **The SDRAM DIMM interface:** This Interface is controlled as described above by the SDRAM Controller and only data are given by the rest of the logic.

A detailed description of the Interfaces of the IP is given in Subsection 7.2.1.

7.2 Functional Implementation

In this section, the functionality of the Queue Manager is described thoroughly. First, the Interfaces of the IP are given in detail, and a description is provided for all the buses and signals. Then, the main data structures of the IP are presented. The commands that are accepted by the Queue Manager are given afterwards as well as a description of their arguments. Finally the Flow Control provisions of the IP are examined.

7.2.1 Interfaces

In the following sections a detailed description of the interfaces of the Queue Manager IP is given.

7.2.1.1 The Interface with the CPU

In table 7-1 the buses and signals of the CPU – Queue Manager Interface are given. Through this interface the CPU can place commands to the Queue Manager and receive the results when necessary.

Table 7-1: The Interface with the CPU

Command [63:0]	Input	This Input bus is used by the CPU to give a command to the Queue Manager for execution.
In_Data [63:0]	Input	If the command requested by the CPU is a Write operation to a 64-bit word in SDRAM, then this bus is used to give the data that are written.
Out_Data[63:0]	Output	If the command requested by the CPU is a Read operation from a 64 bit word in SDRAM, or any other command that requires data to be returned to the CPU, then this bus is used to give to the CPU the data.
CPU_req	Input	When the CPU requests a command from the Queue Manager this signal is used as a flag of the request. Command[63:0] and other buses must be valid.
Com_Acc	Output	When the Queue Manager starts the processing of a CPU command this flag signal is set to notify the CPU that it can prepare a new command since the previous command attributes have been latched in the IP.
Com_ready	Output	When a CPU command has been executed and resulting data are valid on Out_Data[63:0] then this signal is set to notify the CPU.
CPU_lock	Input	For debugging reasons, the CPU can set this signal that orders the Queue Manager to process only CPU commands and stop any enqueueing or dequeueing of cells.

7.2.1.2 The Cell Demux Interface (Incoming cells)

In table 7-2 the buses and signals of the Cell Demux – Queue Manager Interface are given. This interface can be used to give to the Queue Manager the incoming cells that must be enqueued. The Queue Manager is the master of the Interface and decides when the cells can be read in.

Table 7-2 : The Cell Demux Interface (Incoming cells)

Enq_Req	Input	When a cell is available in the 7x64 Input buffer of Cell Demultiplexor then this signal is set to notify the Queue Manager that a cell
---------	-------	---

		needs enqueueing. This signal is the fifo full signal of the buffer.
Cell_DataIn[63:0]	Input	When Enq_Req = 1 then this bus gives the 7 64-bit words of the cell. The first word of the cell is valid on the first cycle the Enq_Req = 1 and contains the Flow ID that the cell belongs to.
Read_En	Output	The Read_En that notifies the Buffer that a cell word has been latched and it must produce the next in the next cycle.
Cell_Read	Output	When a whole cell has been read, this signal is set by the Queue Manager to notify the buffer that it can be loaded with a new cell.

7.2.1.3 The Cell Scheduler – Queue Manager Interface

In table 7-3 the buses and signals of the Cell Demux – Queue Manager Interface are given. This interface can be used by the Scheduling logic to request cells of specific Flow Groups to be dequeued and given to it. If subsequent transmission of the cells fails re-read of the cells can be requested also. A FIFO of free buffers from successfully transmitted cells can be kept in the Cell Scheduler, and its elements can be given to the Queue Manager through this interface, to assist the latter in free-list bypassing.

Table 7-3: The Cell Scheduler – Queue Manager Interface

Deq_Req	Input	When the Cell Scheduler wants to make a request to the Queue Manager, to dequeue a cell from a specific Flow Group, it sets this signal.
Flow_Group[5:0]	Input	When the Cell Scheduler wants to make a request to the Queue Manager, to dequeue a cell from a specific Flow Group, it puts the Flow Group ID on this bus (from 0 to 63)
Cell_DataOut[63:0]	Output	The Queue Manager uses this bus to give the dequeued cell to the Cell Scheduler in 7 64-bit words.
Write_En	Output	The Queue Manager uses this signal to write the in 7 64-bit words into the Cell Scheduler. This signal is used as a write enable in an internal cell buffer kept in the Cell Scheduler.
Flow_Id_Out[15:0]	Output	During the cycle of first word of the outgoing cell, this bus contains the ID of the Flow that the cell belongs to.
Cell_Addr_Out[21:0]	Output	During the transmission of the outgoing cell, this bus contains its buffer address inside the SDRAM memory module. The Cell Scheduler will store this address in case the transmission

		of the dequeued cell fails. In that case it will request from the Queue Manager to read again the cell from the SDRAM, by providing this address.
Read_Cell_Req	Input	If the transmission of a previously cell fails, then the Cell Scheduler uses this signal to request from the Queue Manager to read again the cell from the SDRAM
Cell_Addr_In[21:0]	Input	If the transmission of a previously cell fails, then the Cell Scheduler uses this bus to give the address of the cell that must be read again from the SDRAM The address was given to Cell Scheduler during the first dequeue, through Cell_Addr_Out[21:0] bus.
Free_Cell_Empty	Input	In case a cell is transmitted without error its address in the SDRAM memory is kept in a FIFO that is read by the Queue Manager. This FIFO provides free buffers to the latter for free list bypassing (See Subsection 7.3.7). This signal when set notifies the Queue Manager that this FIFO is empty and no free list bypassing can be done.
Free_Req	Output	If the FIFO of freed buffers inside the Cell Scheduler is almost full, the latter requests with this signal from the Queue Manager to dequeue one free buffer.
Free_Cell_Ptr[21:0]	Input	This bus carries to the Queue Manager the head of the list of free buffers inside the Cell Scheduler (when not empty). The Queue Manager uses this buffer for free list bypassing or to add it to the free list inside the IP, in case the FIFO is almost full ($Free_req = 1$).
Free_Cell_Ld	Output	After reading the Free_Cell_Ptr[21:0], the Queue Manager uses this signal as a read enable of the FIFO to dequeue the buffer and output the next.
Deq_Ack	Output	When a request by the Cell Scheduler ($Deq_Req = 1$, $Read_Cell_Req = 1$, or $Free_Req = 1$) has been accepted by the Queue Manager, the latter sets this signal to 1 for 1 cycle.
Cell_Ready	Output	When a whole cell has been written into the Cell Scheduler (after a Dequeue or a Read Cell request) this signal is set to 1 for 1 cycle.

7.2.1.4 The SDRAM – Queue Manager Interface

The last interface of the Queue Manager presented here, is the only one that is external to the FPGA that hosts it. It is the pin interface of the Queue Manager with

the SDRAM DIMM memory module. Through this interface, the IP makes its memory accesses, in order to enqueue or dequeued cells, or to enumerate the data structures kept there. The control signals and address of this Memory interface is set by the SDRAM Controller sub-block of the IP while the data are given or taken directly by the Queue Manager Datapath. The interface is the SDARM DIMM standard and any commercial SDRAM can be used. Table 7-4 gives the control and data signals-busses of this interface.

Table 7-4: The SDRAM – Queue Manager Interface

MemData[63:0]	InOut	The bi-directional Data Bus of the Interface, used to read/write 64-bit words from/to the Memory
SDRAMaddr[11:0]	Output	The address bus, used by the Queue Manager to give the row address and the column address of an access
BnkEn[1:0]	Output	The signals that select the internal banks of the SDRAM chips on the DIMM
ClkEn[1:0]	Output	These Clock Enable signals activate or deactivate the DIMM
Cs[3:0]	Output	These Chip Select signals select or deselect half of the SDRAM chips on the DIMM.
We_	Output	Write Enable signal. When low, a write access is performed.
RAS_	Output	Row Access Strobe, low when row address is given.
CAS_	Output	Column Access Strobe, low when column address is given.
DQM[7:0]	Output	Mask bits for each of the 8 bytes of the 64 bits words. Not used by Queue Manager.

7.2.2 Flows

The Flow is the main data structure that is supported by the Queue Manager. Its meaning is that of a flow of connection traffic that is served (its cells are enqueued or dequeued) by the Queue Manager. The Queue Manager can support up to 64K flows of traffic. To distinguish between them its Flow has a specific ID (Flow ID) number that is reserved by the CPU during connection-flow set up. Since the maximum number of supported flows is 64K the Flow ID is a 16-bit binary number. Information for each of the flow is kept in Flow Records. These Records have always reserved space inside the SDRAM memory. The Data Structure that implements the Flow is that of the unidirectional list seen in figure. The Record of each flow maintains information for that list like, the head pointer (address of the first element's buffer), tail pointer (address of the last element's buffer), counter (number of cells in the list), status bits (if the flow ID is used by a connection, or if the Flow ID is active) and Quality of Service parameters. Note that a flow is considered used when the CPU has reserved it for a connection during connection setup and active when it has a cell

stored inside the SDRAM. A detailed description of the Flow Record is given in Subsection 7.3.1.

When a cell that belongs to a certain flow (this is determined by the Flow ID that it carries in its header) arrives at the Queue Manager, it is enqueued in the list. The Queue Manager does this, by writing the cell into the free buffer that is always at the end of the list, and setting the next pointer of this buffer to point to a new empty buffer (There must always be an empty pointer at the end of the list. This is due to reasons explained in Subsection 7.3.8). Empty buffers for this purpose are taken from a list (Free List) that is kept by the Queue Manager. The Tail Pointer of the Flow Record is set to point to that empty buffer also and the counter is increased.

When a cell from that flow must be dequeued, and send to the Cell Scheduler for transmission over the Cell Bus, then the Queue Manager reads the cell that is the head of the list and sets the head pointer to point at the next buffer. If the transmission of the cell succeeds then the buffer that was read can be put to the Free list or reused as the empty tail buffer of an enqueue operation (Free List bypassing. See Section 7.3.7). Figure 7-2 depicts the Flow structure, the enqueue operation and the dequeue operation.

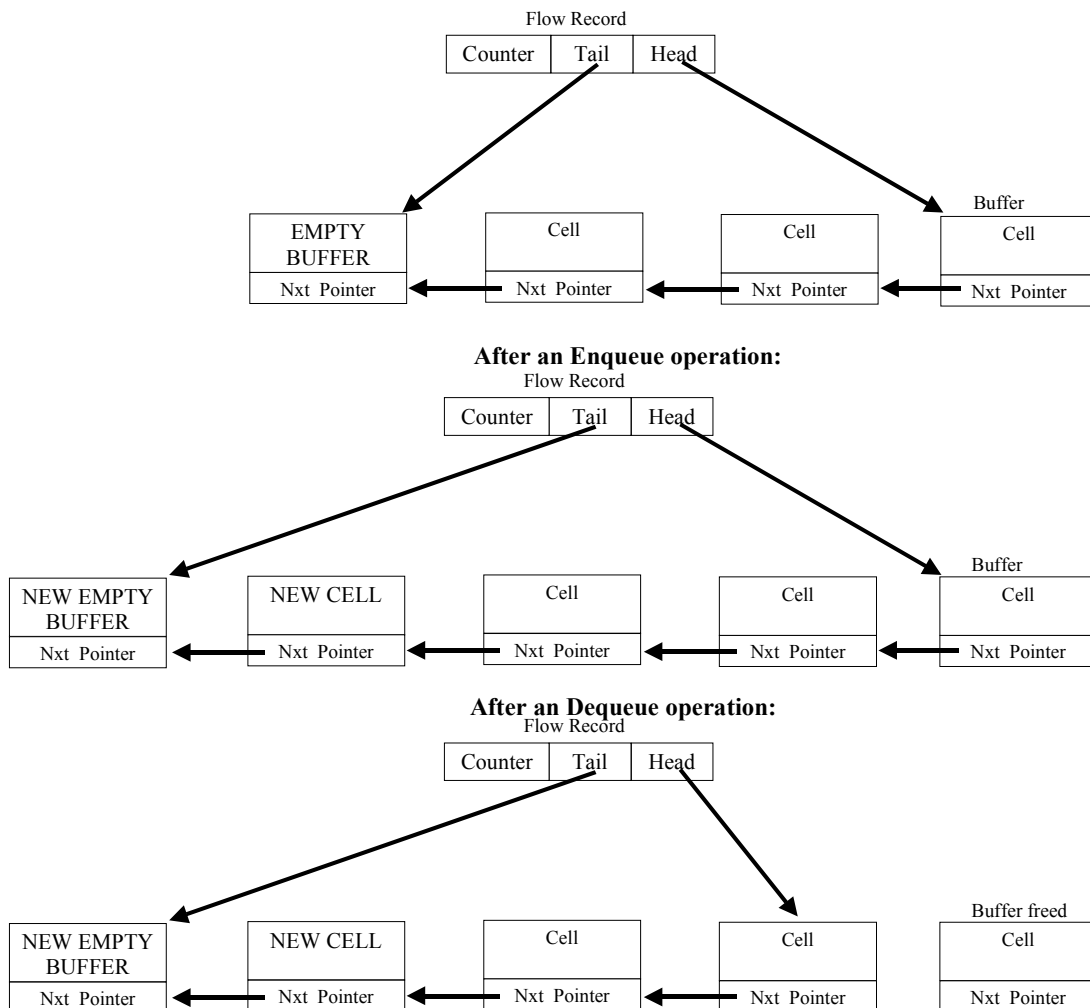


Figure 7-2: The Flow Structure and changes on it after an enqueue and a dequeue operation.

7.2.3 Flow Groups

The 64K flows supported by the Queue Manager are organized in a higher level to Flow Groups. All the flows that are used by connections must belong to a flow group. The number of Flow Groups supported by the Queue Manager is 64. The reason of this grouping is to assist the Scheduling block in the schedule of outgoing ABR traffic. Per flow scheduling according to QoS parameters is impossible for a very large number of flows. For that reason flows are organized in flow groups by the Queue Manager, each Flow Group has its own QoS parameters kept in the Cell Scheduler and the latter has to schedule traffic by providing bandwidth to these 64 Flow Groups. It does this by requesting from the Queue Manager to dequeue cells from specific Flow Groups. The Flows in each Flow Group each get equal amount of bandwidth, since the Queue Manager serves them in a round robin manner.

The data structure that implements a Flow Group is that of a cyclic list depicted in Figure 7-3. Heads and tails of this list (there are 64 of them) are stored in the Flow Group Memory. Only active flows exist in these lists. If a flow after dequeuing its last cell in SDRAM memory becomes empty then it is removed from the cyclic list. When a new cell arrives for that flow, it is then reintroduced in the list by being placed as the tail of the list.

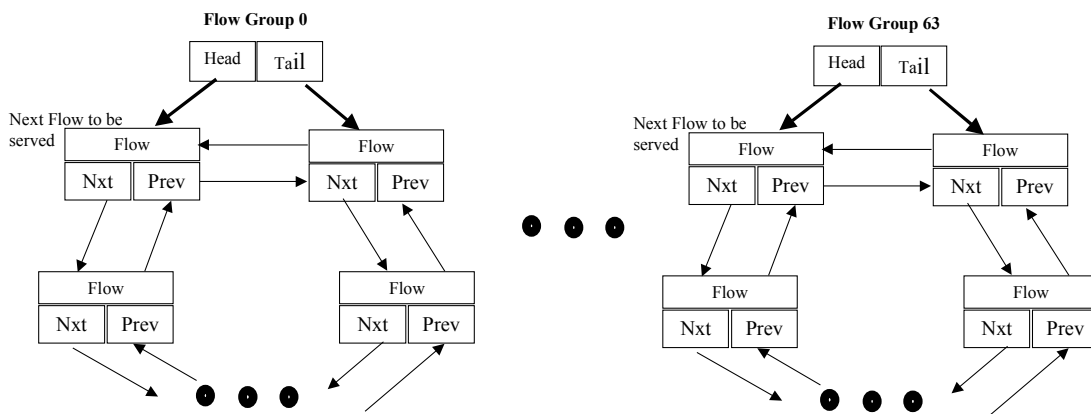


Figure 7-3: The 64 Flow Group cyclic lists

If the Cell Scheduler requests a dequeue operation for a specific Flow Group, the Flow that is head of the cyclic list will be selected to dequeue a cell and the Tail Pointer will be set to point to that Flow. Then the Head pointer of the cyclic list will point to the next Flow.

The cyclic list is unidirectional since the connection represented by a Flow may close. In that case the flow must be removed from the list while the list must remain connected in a $O(1)$ operation. Pointers to the next and the previous flow for a specific flow are kept in that flow's Record.

7.2.4 Queue Manager Commands

The Queue Manager can accept a range of commands and execute them using one State Machine for each one of them, to control data structures, registers and memory

for each one of them. The most important commands are the Enqueue and Dequeue Commands, issued by the Cell Demultiplexor and Cell Scheduler respectively. Their importance lies in that the number of cycles that spent for their execution defines the Queuing bandwidth of the Queue Manager. Table 7-5 lists the available commands of the Queue Manager, their arguments, their return data and the number of cycles needed for their execution.

Table 7-5: Table of all the Queue Manager commands

NAME	Issued by	Args	Return Data	Clks	Description
Read	CPU	Address	Mem Data	5	Read a 64-bit word from the SDRAM Memory.
Write	CPU	Address, Data		5	Write a 64-bit word to the SDRAM Memory.
OpenFl	CPU	FlowID, FGID, Hwmark, LWmark		10	Initialize a Flow at connection set-up. Reserves a Flow ID, and assigns it to a Flow Group. Sets the Flow Control parameters(Hwmark, Lwmark)
CloseFl	CPU	FlowID		20	Closes a flow at connection shutdown. Releases the Flow ID removes it from Flow Group cyclic list and adds all used buffers to free list.
ReadCnt	CPU	FlowID	Counter	5	Read the counter of buffered cells of a Flow from its record in the SDRAM Memory.
Enqueue	Cell Demux	FlowID, Cell		20, 40	Enqueue of an incoming cell to its respective Flow ID Queue.
Dequeue	Cell Sched	FGID	Address, Cell	20, 40	Dequeue a cell from the Flow that is Head to the cyclic list of the given Flow Group ID. Send it to Cell Scheduler along with its buffer address.
RdCell	Cell Sched	Address		12	Read the contents of the buffer with the given address.
Free	Cell Sched	Address		10	Add the buffer with the given address to the Free List.
ChParam	CPU	FlowID, Hwmark, LWmark		10	Change the Flow Control parameters of the Flow ID to the given.

The Enqueue requires 20 clock cycles to be executed, except in the special case that the flow was inactive. In that case the Flow must enter into the cyclic list of its respective Flow Group. The records of its previous and next flows in that list must be enumerated adding another 20 cycles in the execution cycles. Thus the number increases to 40.

The Dequeue also requires 20 clock cycles to be executed, except in the special case that the flow becomes inactive (the cell dequeued was the last in SDRAM memory). In that case the Flow must be removed from the cyclic list of its respective Flow Group. The records of its previous and next flows in that list must be enumerated adding another 20 cycles in the execution cycles. Thus the number increases to 40.

7.2.5 EFCI and RM marking [21]

One of the features of the Queue Manager IP is the provision of generic ATM FORUM Flow Control mechanisms. There are two ways that misbehaved flow can be ordered to drop its traffic flow, according to the ATM FORUM specifications.

The first is EFCI (Explicit Forward Congestion Indication). EFCI is a bit included in the header of each ATM cell. When congestion exists or is about to exist in a network node, the latter can set this bit to one, ordering thus the traffic destination to decrease its request of data.

The second is RM (Resource Management) cell marking. RM cells are regularly sent by the source of ABR traffic to the destination node. The destination U turns these cells back to the destination. These cells contain information about traffic resources of the intermediate nodes. Two bits inside this cell, CI (Congestion Indication) and NI (No Increase) can be set to drop or stop the increase of the traffic flow respectively.

The Queue Manager uses both of these mechanisms to instruct a misbehaved flow to decrease its cell rate. A misbehaved rate is determined by a Hi Watermark – Low Watermark mechanism. At the time of initialization of a Flow Record during connection set-up, the CPU set the highest number of cells (Hi Watermark) that the Flow is allowed to use in the SDRAM memory. This Hi Watermark is stored inside the Flow Record. When the Flow Counter during an enqueue operation on that Flow surpasses this Hi Watermark, a bit inside the Flow Record (Mark bit) is set, and the cells that are dequeued from that time and on are EFCI marked. If the cells are RM their CI, NI bits are also marked. Cells are marked at the output in order for the congestion information to reach its target faster. This saves us the time of waiting the number of unmarked cells to leave the Flow Queue.

The cells will stop being marked when the Flow Counter becomes less that Low Water Mark that is also set during the time of Flow initialization. Resetting the Mark bit in the Flow Record does this.

The ChParam command issued by the CPU and described in subsection 7.2.4 changes the Hi Watermark – Low Watermark fields in the Flow Record, thus dynamically affecting the Flow Control mechanism. This means that a well-behaved flow, after the execution of this command can become misbehaved and the marking of its dequeued cells will commence. A detailed description of the fields in the Flow Record, relative to the flow control mechanism of the Queue Manager is given in subsection 7.3.1.

7.3 Design Implementation

7.3.1 Flow Record Format

In figure 7-4 the detailed description of the Flow Record is given. There are 64K records of like this, one for each of the 64K Flows. Each record is 2 words of 64 bits long. This is 2 SDRAM memory words. They are stored inside the SDRAM memory in 2 word alignment.

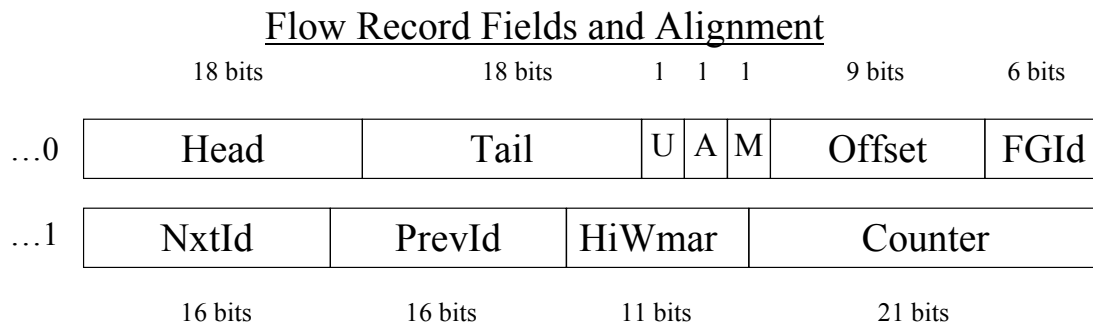


Figure 7-4: Flow Record Fields and alignment

Table 7-6 gives a description of each field in the Flow Record.

Table 7-6 : Flow Record Field bits and description

Field	Bits	Description
Head	22	Head Pointer: It contains the address of the head of the queue of the respective flow.
Tail	22	Tail Pointer: It contains the address of the tail of the queue of the respective flow.
Used	1	Used Flow: When 1 this Flow ID is used by a connection. All the record fields are valid.
Active	1	Active Flow: When 1 this Flow ID has not empty queue. Some cells belonging to the flow are stored in the SDRAM memory.
Mark	1	Mark Bit: The cell count for the respective queue has surpassed the Hi Watermark. Dequeued cells are RM, EFCI marked. Will reset to 0 when Cell Counter drops below Hi Watermark - Off
HiWmH	6	Hi Watermark Most Significant bits: The Most significant bits of the Hi Watermark.
HiWmL	12	Hi Watermark Least Significant bits: The Least significant bits of the Hi Watermark.

FGID	6	Flow Group ID: The Flow Group that this Flow is assigned.
Off	5	Offset: Hi Watermark – Off is equal to the Low Watermark
NextID	16	Next Flow ID: The next flow in the cyclic list of this flow's Flow Group. If Active = 0 this field is invalid.
PrevID	16	Previous Flow ID: The previous flow in the cyclic list of this flow's Flow Group. If Active = 0 this field is invalid.
Counter	20	The number of cells of the flow inside the SDRAM memory.

An effort was made during the design so that the Flow record would be 2 words long. If it was 3 words long it would mess the 2 word alignment of the SDRAM memory. If 4 words it would require 2 extra cycles to access it, dropping the queueing bandwidth. Buffer pointers such as Head and Tail are buffer aligned in a 256 MByte Memory organized in 2^{25} words of 64 bits. A buffer is 8×64 bit words, thus a buffer aligned pointer is 22 bits long.

Notice that there is no field for the Flow ID of the Flow that the Record belongs to. This is done because each Flow Records is stored in 2 word whose address is that of the Flow ID. An additional least significant bit is used in the addressing of the record to distinguish between the 2 words of the record. In that organization, we put the records at the beginning of the SDRAM memory (See Sub-Section 7.3.4) and save field space in the record.

7.3.2 Flow Group Record Format

All the Flow Group information is stored in a 64×33 memory, the Flow Group memory. Each of the 64 words of this memory stores the Flow Group record with ID that of its address and keeps all the data necessary for the maintenance of the Flow Group cyclic list. Figure 7-5 depicts the Flow Group memory and the Flow Group records.

	16 bits	16 bits	1
0	Head Flow	Tail Flow	A
1	Head Flow	Tail Flow	A
2	Head Flow	Tail Flow	A
3	Head Flow	Tail Flow	A

62	Head Flow	Tail Flow	A
63	Head Flow	Tail Flow	A

Figure 7-5: Flow Group memory organization and Flow Group records.

Table 7-7 gives the description of the Flow Group Record fields.

Table 7-7: Flow Record Field description

Field	Bits	Description
Head Flow	16	Head Flow ID: It contains the Flow ID of the Flow that will give the next dequeued cell, when a dequeuing operation is requested for the respective Flow Group.
Tail	16	Tail Flow ID: It contains the Flow ID of the Flow that was served in the previous dequeuing operation for that flow group.
Active	1	Active Flow: When 1 this Flow Group has Active Flows, else no Flow assigned to this Flow Group has any cells inside the SDRAM memory

The contents of the Flow Group memory are visible to the Cell Scheduler. The latter needs to know which of the Flow Groups are active so that it will keep them in the scheduling loop. Request for dequeuing of an inactive Flow Group would cause an error.

7.3.3 Cell Format and Alignment

Each cell is stored inside the SDRAM memory in 8x64 buffers. The 7 first words (7x8=56) store the ABR cell plus the internal Switch header (Cell Bus, Tandem Routing headers in the DIPOLo environment). The last word of the buffer stores the next pointer. This points to the buffer that stores the next cell of the flow queue. Each

buffer is 8 word aligned. Figure 7-6 depicts the Cell Format and Alignment in the SDRAM Memory.

...000	Cell 0
...001	Cell 1
...010	Cell 2
...011	Cell 3
...100	Cell 4
...101	Cell 5
...110	Cell 6
...111	Next Ptr

Figure 7-6: Cell Format and alignment

7.3.4 SDRAM Memory Organization

The SDRAM DIMM module that supports the memory needs of the Queue Manager has a capacity of 256 Mbytes. Inside this memory space the Flow Records for each of the 64K flows is stored. The rest is divided in cell buffers that are dynamically allocated to the incoming traffic. The Flow Records are statically allocated. This means that all of the 64K Records are present, even if they are not used.

The CPU is capable of initializing the SDRAM memory by using the Write command. Still, there is an FSM inside the State Machine that after reset can initialize the SDRAM memory. The contents of the SDRAM memory after proper initialization are shown in Figure 7-7.

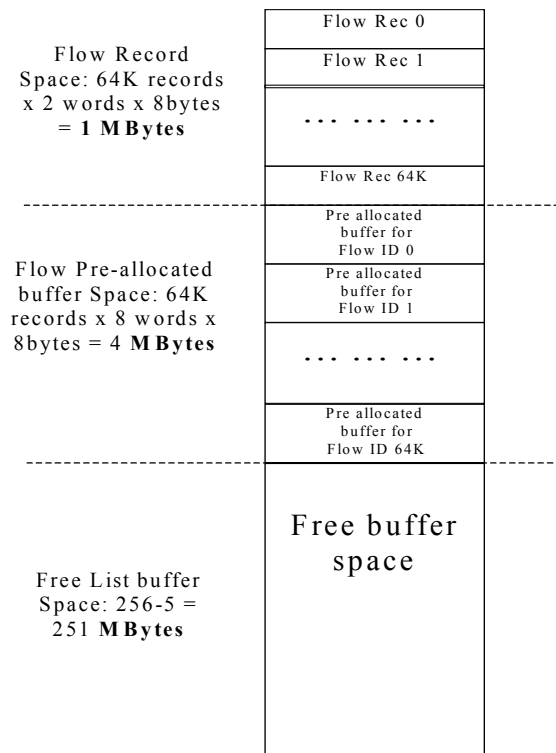


Figure 7-7: SDRAM Memory-space division and organization

Since the memory is 256 Mbytes = 2^{28} bytes, organized in 64-bit (8 bytes) words, the number of word addresses in the memory is 2^{25} . Pointers to a buffer that is 8 words and aligned need to be 22 bits.

As described in subsection 7.3.1, in order not to have a field for the Flow ID of each Record, the ID is used as a pointer to the position of its Flow Record in the Memory and the Records are placed in the beginning of the Memory. Thus, the address of the first word of the Flow Record of Flow ID 0b1111000111110001 (Flow Id is 16 bits) is 0b{00000000,1111000111110001,0}, while the address of the second word is 0b{00000000,1111000111110001,1}. Since the Queue manager supports 2^{16} Flows (64K) then the Flow Record use the first 2^{16} Flows * 2 words/Flow = 2^{17} words = 2^{20} bytes = 1 MByte of Memory.

Since each Flow Queue must always have an empty buffer (even if not used, see subsection 7.3.8) due to buffer pre-allocation, one is given to each during initialization. These buffers are positioned in the SDRAM memory, after the Flow Records. After some time of system operation these buffers are used for cell storing, but others take their place as empty buffers. Thus the number of empty pre-allocated buffers is constant and equal to 2^{16} (one for each Flow). So buffer pre-allocation uses up another 2^{16} buffers * 8 words/buffer = 2^{19} words = 2^{22} bytes = 4 MBytes.

The remaining $256 - (4+1) = 251$ MBytes can be used for cell buffering. During Memory Initialization, this free space is organized in a large FIFO unidirectional list of free buffers, the Free List. This Free List is used to provide the enqueue operation with free buffers for storing and accept the freed buffers after the dequeue operations.

7.3.5 State Machines

As described in section 7.1, the State Machine sub-block of the Queue Manager contains the State Machines that execute the commands requested by the other blocks. Figure 7-8 shows the internal hierarchy of the FSMs inside this sub-block.

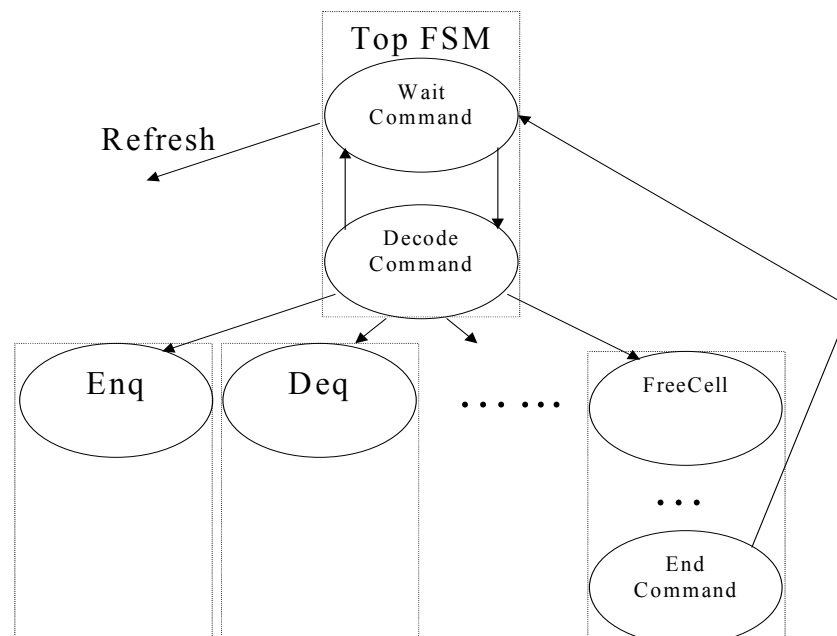


Figure 7-8: State Machine Top Level Diagram

These FSMs set the control signals of the Datapath of the Queue Manager, they also set and accept the control signals of its Interfaces and the make requests for SDRAM memory accesses to the SDRAM controller block. There is one FSM for each one of the commands of the Queue Manager IP. At the top of the hierarchy of these FSMs lies the TOP FSM. The purposes of this FSM are:

- **Initialize the Datapath Registers after System Reset:** Some registers like (Free List Head) need to be initialized in a specific value after reset. The TOP Fsm sets the control signals to do that.
- **Accept the Command Requests by other blocks:** The TOP FSM remains idle, and pools the request signals from the command issuing blocks.
- **Arbitrate Queue Manager Command issue:** There is the case that more than several blocks can request one command, at the same time. The TOP Fsm arbitrates which one will use the Queue Manager, according to some priorities. Initiate a command FSM: When a command is to be executed, the Top Fsm sets the control signal that initiates its respective State Machine. It waits for the control signal from that State Machine that notifies that the execution has finished. It also acknowledges the command acceptance to the issuing block.
- **Request for a SDRAM Memory refresh:** SDRAM Memories must be refreshed periodically. The TOP Fsm has an internal counter that when it reaches zero, issues SDRAM Memory refresh command to the SDRAM controller.

Table 7-8: Queue Manager Command Priorities

PRIORITY	NAME	Issued by	Comments
1 (highest)	Refresh	State Machine	Needs to be executed periodically or else data in SDRAM may be lost.
2	Write, Read, OpenFl, CloseFl, ReadCnt, ChParam	CPU	CPU Commands have the same priority for the reason, that only one can be requested at any time. They have higher priority than the ones issued by Cell Demux, Cell Scheduler because they are sort, they configure and control the Queue Manager operation and are necessary for debugging purposes.
3	Enqueue,	Cell Demux	Enqueue command (issued by Cell Demux) has equal priority with the commands issued by Cell Scheduler. When there is a contention between Enqueue and one Cell Scheduler command, an alternative arbitration is made by TOP FSM. If most recent command execution was Enqueue, the Sched command will be selected, else if most recent was Sched Command then Enqueue executes
	Dequeue, RdCell, Free	Cell Scheduler	

Table 7-8 gives the prioritization of Commands that is taken into account when there is a contention that must be resolved by the TOP Fsm arbitration. The command priority is given from highest to lowest. Commands in the same row have the same priority.

Figure 7-9 shows a simplified state diagram of the FSM that executes the Enqueue Command. This FSM along with the Dequeue FSM are the most complex in the State Machine Sub-block. The states where an access to the SDRAM is requested to the SDRAM controller are shown with big bubbles that contain the type of access. The maximum number of states/cycles needed is 40. This is the case when the flow was previously inactive. In that case the Flow must enter into the cyclic list of its respective Flow Group. The records of its previous and next flows in that list must be enumerated adding another 20 cycles in the execution cycles. Thus the number increases to 40. Some 5 cycles can be saved when Free-List bypassing is being implemented (See Sub-section 7.3.7). In that case an access to get a new free buffer pointer is avoided.

The SDRAM memory accesses are being made in such an order that data dependencies are avoided (the result of a read operation is used as the address of the next read or write access immediately), while the SDRAM is used continuously.

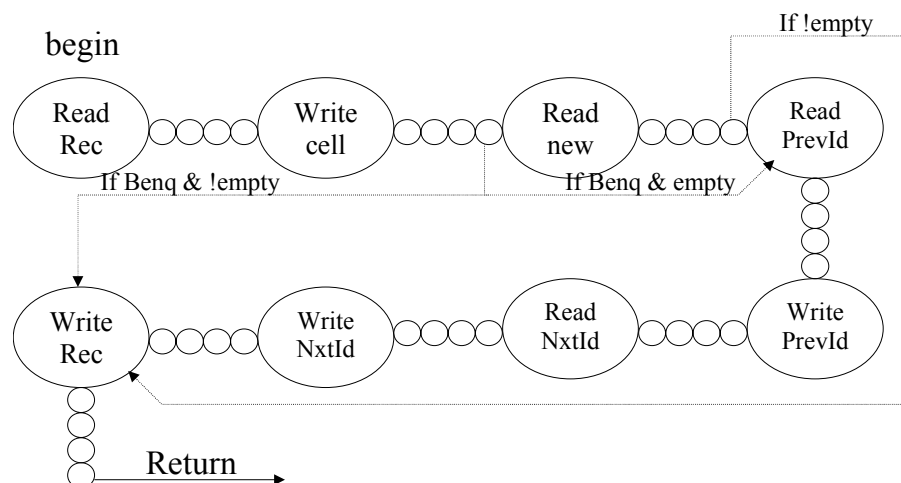


Figure 7-9: Enqueue command FSM bubble diagram

7.3.6 The SDRAM Controller

The SDRAM controller is the sub-block of the Queue Manager that is responsible for controlling the SDRAM memory. It accepts simple burst commands from the State Machines of the Queue Manager and issues these to the memory. The addresses needed are received by the Queue Manager datapath and are used to set Chip Select Signals, Bank Enable signals, as well as to provide the Row and Column addresses for the burst access. The mode register that configures the operation of the SDRAM DIMM is also written during initialization by the sub-block.

The Queue Manager uses the SDRAM Controller for the following operations:

- **Load the Mode Register:** This is done after Reset of the system
- **Auto Refresh:** The TOP Fsm periodically requests the refreshing of one SDRAM row.
- **Read burst of 1 64-bit word:** Read access to half of a Flow Record.

- **Read burst of 2 64-bit words:** Read access to a Flow Record.
- **Read burst of 8 64-bit words:** Read access to Cell Buffer plus its next pointer.
- **Write burst of 1 64-bit word:** Write access to half of a Flow Record.
- **Write burst of 2 64-bit words:** Write access to a Flow Record.
- **Write burst of 8 64-bit words:** Write access to Cell Buffer plus its next pointer.

Figure 7-10 gives the State machine diagram of the SDRAM Controller block:

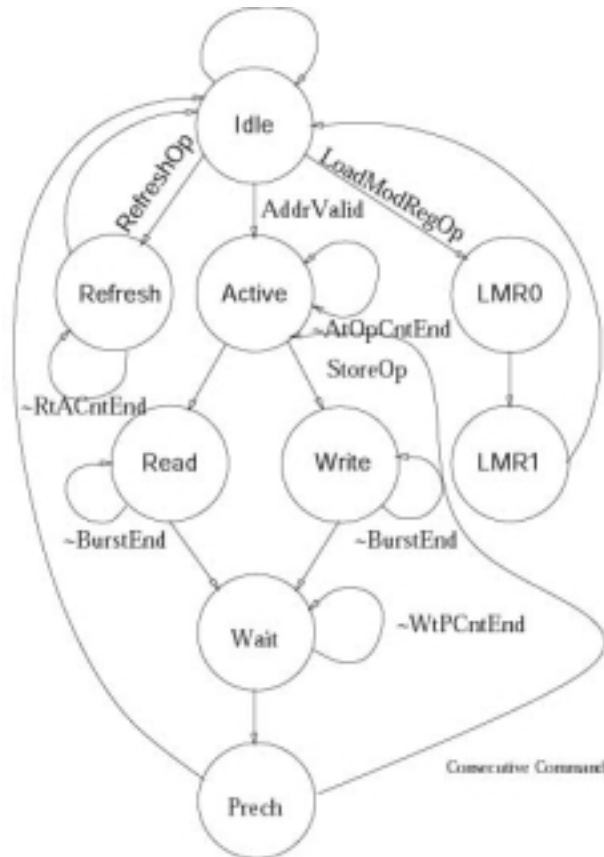


Figure 7-10 : State Machine Diagram of the SDRAM controller

The Controller is in idle state until, a read or a write command is requested by the State Machines. Then it activate the accessed row to be accessed and makes the burst access (stays in Read, Write states for as many cycles as the burst size), waits one cycle and pre-charges the used raw. If a consecutive command has already been requested, the state machine enters to the Active state again, or else in moves to the idle state.

The SDRAM controller doesn't handle the data to/from the SDRAM DIMM. This is done by the Queue Manager datapath, which places them to, or reads them from, the SDRAM 64-bit data bus, in sync with the SDARM Controller State.

7.3.7 Free List Bypassing [13]

The Free-List Bypassing is a Queueing technique implemented by the Queue Manager IP that avoids an additional access to SDRAM memory during an Enqueue or a Dequeue Operation. This decreases the number of cycles needed for both to finish execution and thus increases the Enqueueing/Dequeueing Bandwidth.

When an Enqueue operation is being done a new cell Buffer must be given to the Queue that receives the cell. The Cell Buffer that is pointed by the Free List Head pointer is selected for this purpose but now the Free List Head pointer must be set to point to the next free cell buffer of the Free List. Reading the Next Pointer of the cell buffer taken from the Free List does this. This Next Pointer lies in the External SDRAM Memory and that read access will cost an additional 5 cycles to the Enqueue process.

Moreover, when a Dequeue operation is executed and the dequeued cell is transmitted correctly, its cell buffer is no longer used and must be added to the Free List. This is done in the Queue Manager with the execution of the Free Command that enumerates the Free List Head pointer with the address of the newly added cell and write in the Next Pointer of the latter the address of the previous Free List Head cell. This Write access is done to the SDRAM Memory and cost 5 cycles.

Free list bypassing avoids the cost of the extra 5 cycles spent on Enqueueing and on Dequeueing, as described above. Instead of placing the newly freed cell buffer in the Free List (5 cycles of Free Command), the pointer is placed in a SRAM FIFO inside the Cell Scheduler. When a subsequent Enqueue operation begins, the cell new cell buffer requested, is taken from that SRAM FIFO, thus the extra access to the Free List is also avoided.

In total the Free List Bypassing decreases the number of cycles for enqueueing and dequeueing of a cell by 10 cycles. Since Enqueueing as dequeueing takes 20 or 40 cycles (See subsection 7.3.5), the technique improves the Queue Manager Performance by a factor of $10/(20+20) = 1/4 = 25\%$ or $10/(40+40) = 1/8 = 12,5\%$.

Figure 7-11 depicts the Free-List Bypassing technique

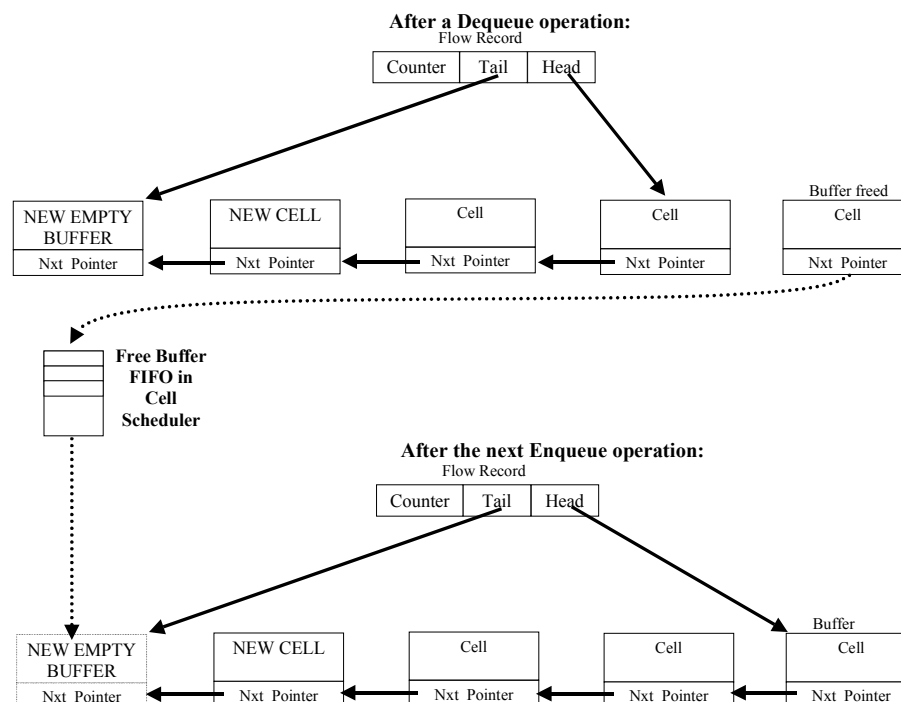


Figure 7-11: Free List Bypassing implementation in the Queue Manager IP

7.3.8 Cell Buffer Pre-allocation [10]

The Cell Buffer Pre-allocation is a technique implemented in the Queue Manager that with the cost of one unused cell buffer per supported Flow, decreases the number of cycles for an enqueue operation by 5.

As described in previous sections, each one of the 64K queues of the Queue Manager has an empty cell buffer as its tail element. This is true even for queues that are empty (inactive). In this case both the Head and the Tail pointer field of the respective Flow Record point to that empty cell buffer (this buffer is called the pre-allocated buffer). If this buffer wasn't pre-allocated and the tail pointer of the Flow Record pointed to the last used cell buffer, then during enqueueing, the newly arrived would have to be written on a new empty cell, taken from the free-list (or the Free-List Bypassing mechanism), and the next pointer of the previously last cell should be set to point to the newly arrived cell. These are two memory accesses in different buffers in memory. Since these two buffers can be in different SDRAM lines, two separate bursts would be necessary.

Instead, with Cell Buffer Pre-allocation, the cell is written to the pre-allocated empty cell buffer that is tail element and the next pointer of the same buffer is set to point to the new pre-allocated buffer taken from the free-list (or the Free-List By-passing mechanism). Only one buffer is accessed in this way, and only one write burst is needed. The extra pointers write access in the non pre-allocating implementation that is avoided costs 5 cycles.

In total, the Free List Bypassing decreases the number of cycles of enqueueing and dequeueing of a cell by 5 cycles. Since Enqueueing as well as dequeueing take 20 or 40 cycles (See sub-section 7.3.5), the technique improves the Queue Manager Performance by a factor of $5/(20+20) = 1/8 = 12,5\%$ or $5/(40+40) = 1/16 = 6,25\%$. Memory cost of this technique is 4 Mbytes which is $4/256 = 1/64 = 1,56\%$ percent of the total memory.

7.4 Timing Issues

In this subsection the bandwidth capabilities of the Queue Manager IP are examined in conjunction with the input/output bandwidth requirements of the ABR Server Cards. The timing parameters assumed and the timing examples that are given prove that the implementation of Free-List Bypassing, along with Cell Buffer pre-allocation are essential if the Queue Manager is to succeed in complying with the Switch requirements, while using only one SDRAM memory module for all necessary memory accesses.

7.4.1 UTOPIA clock Vs Queue Manager (ABRSU) clock.

In order for the Queuing Architecture implemented to be able to balance the incoming and outgoing capabilities of the Switch (DIPOLLO), it must be able both to enqueue and dequeue a cell during a cell time. A cell time is the time needed for a cell to arrive or to depart in full in or out of the UTOPIA interfaces. These interfaces are slaves to the switching hardware. Since enqueue and dequeue operations can't be performed in any parallel way and they both need the same number of cycles, the Queue Manager

(and the ABRSU) must be use internal speed up in relation to the UTOPIA Interfaces clock.

A cell needs 28 UTOPIA clock cycles to enter the ABRSU (16-bit interface). On the other hand, both one enqueue and one dequeue operation need:

80 (Worst case) or 40 (Normal Case) Queue Manager clock cycles.

This means that the clock speed up needed for both cases is:

Worst case:

1 Cell arrival time = 1 Enq time + 1 Deq time=>

$T_{clk_utopia} * 28 \text{ utopia_cycles} = T_{clk_qm} * (40+40) \text{ qm_cycles} \Rightarrow$

$T_{clk_utopia} / T_{clk_qm} = 80 / 28 = 2,8 \Rightarrow$

Speed_up_worst = 2,8

Normal case:

1 Cell arrival time = 1 Enq time + 1 Deq time=>

$T_{clk_utopia} * 28 \text{ utopia_cycles} = T_{clk_qm} * (20+20) \text{ qm_cycles} \Rightarrow$

$T_{clk_utopia} / T_{clk_qm} = 40 / 28 = 1,4 \Rightarrow$

Speed_up_normal = 1,4

The FPGA selected during the design face to host the Queue Manager can reach clock frequencies of up to 50 MHz. Since the Cubit Clocks are used in the DIPOLLO Switch reach 25 MHz of Frequency a maximum speed up value of 2 can be selected. In that case

QM_Clk = 2 * UTOPIA_Clk.

Since:

Speed_up_worst > Speed_up_sel = 2 > Speed_up_normal

The next subsections will show that the Queue Manager is the bottleneck of the system when both enqueue and dequeue take the worst case number of cycles and that the UTOPIA interfaces are the bottleneck system when both enqueue and dequeue take the normal number of cycles.

7.4.2 Worst case of Enqueue and Dequeue

Figure 7-12 gives the timing diagram cells entering and leaving the ABRSU when Enqueue and dequeue commands need 40 cycles each to execute.

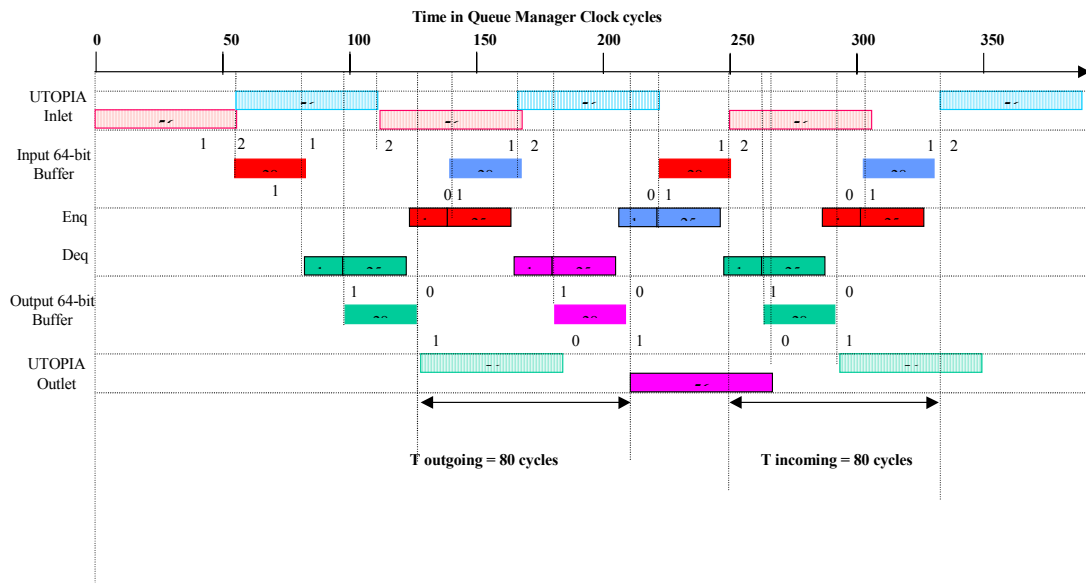


Figure 7-12: Worst case Enqueue-Dequeue timing Diagram

The following parameters have been assumed on this diagram:

- **2 Cell FIFOs in the UTOPIA interfaces. (inlet & outlet) :** The FIFOs inside the UTOPIA interfaces are assumed to be 2 cells long for simplicity. They are 8 cells long - **Single buffering on each Cell Buffer**
- **Speed-up = 2:** The ABRSU speed up is 2. This means that $QM_Clk = 2 * UTOPIA_Clk$
- **Full throttle Traffic:** Incoming cells arrive consecutively, Cell Scheduler requests outgoing cells consecutively.
- **Mixed Enq & Deq Commands:** Enqueue and dequeue commands are executed alternately.
- **40 cycle Enqueue and Dequeue Commands:** Each Enqueue and Dequeue Command takes 40 cycle.
- **15 Cycles access of Input/Output 64-bit Buffers:** The Enqueue command uses its first 15 cycles to empty the Input 64-bit Buffer. The buffer can begin to fill immediately. The Dequeue command uses its first 15 cycles to fill the Output 64-bit Buffer. The buffer can begin to empty immediately.
- **56 QM_clk cycles for UTOPIA I/Fs:** A cell takes 56 QM clock cycles to enter since speed up is 2.

In the diagram, we can see that after 2,3 cell times the UTOPIA inlet is unable to receive the cells consecutively and the UTOPIA outlet to send the cells consecutively. Instead of 56 QM_clk cycles, cells need 80 cycles to enter or leave the ABRSU. This means that the 80 cycles needed by the Queue Manager to Enqueue and dequeue a cell is the bottleneck of the system and sets the cell time to 80 QM_clk cycles.

7.4.3 Normal case of Enqueue and Dequeue

Figure 7-13 gives the timing diagram cells entering and leaving the ABRSU when Enqueue and dequeue commands need 20 cycles each to execute.

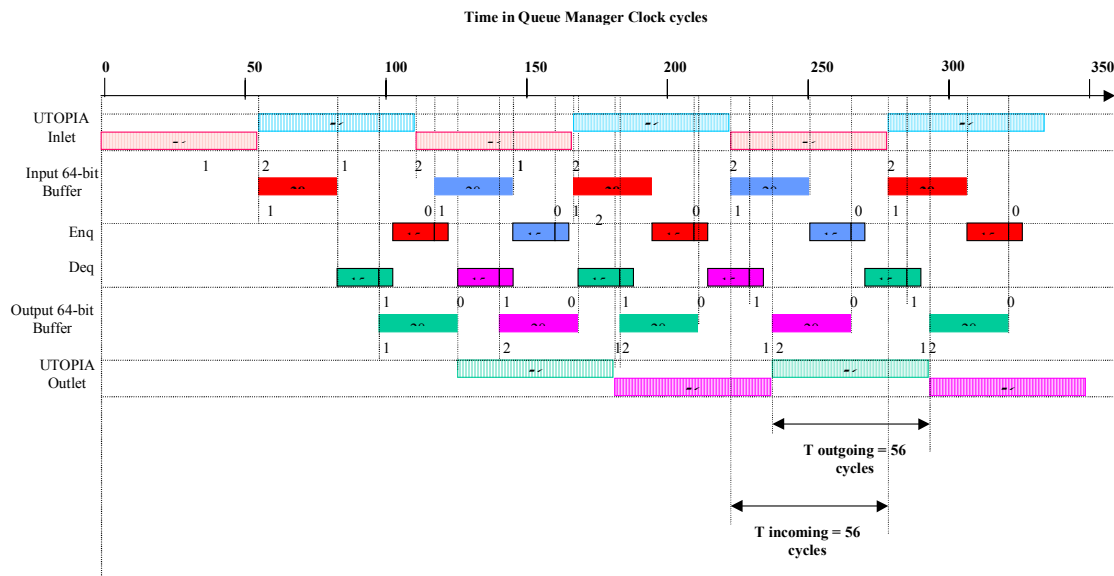


Figure 7-13: Normal case Enqueue-Dequeue timing diagram

The following parameters have been assumed on this diagram:

- **2 Cell FIFOs in the UTOPIA interfaces. (Inlet & outlet) :** The FIFOs inside the UTOPIA interfaces are assumed to be 2 cells long for simplicity. They are 8 cells long - **Single buffering on each Cell Buffer**
- **Speed-up = 2:** The ABRSU speed up is 2. This means that $QM_Clk = 2 * UTOPIA_Clk$
- **Full throttle Traffic:** Incoming cells arrive consecutively, Cell Scheduler requests outgoing cells consecutively.
- **Mixed Enq & Deq Commands:** Enqueue and dequeue commands are executed alternately.
- **20 cycle Enqueue and Dequeue Commands:** Each Enqueue and Dequeue Command takes 40 cycles to complete.
- **15 Cycles access of Input/Output 64-bit Buffers:** The Enqueue command uses its first 15 cycles to empty the Input 64-bit Buffer. The buffer can begin to fill immediately. The Dequeue command uses its first 15 cycles to fill the Output 64-bit Buffer. The buffer can begin to empty immediately.
- **56 QM_clk cycles for UTOPIA I/Fs:** A cell takes 56 QM clock cycles to enter since speed up is 2.

In the diagram, we can see that after 5 cell times the UTOPIA inlet is able to receive the cells consecutively and the UTOPIA outlet to send the cells consecutively. Cells need 56 cycles to enter or leave the ABRSU. This means that the 40 cycles needed by the Queue Manager to Enqueue and dequeue a cell is not the bottleneck of the system rather than the 56 cycles cell time of the UTOPIA interfaces.

7.4.4 Synthesis Results, Free-list bypassing and Cell Buffer Pre-allocation contribution

After synthesis of the Verilog HDL files, that described the ABRSU blocks with the MaxPlusII synthesis tool, the following results were received.

- **ABRSU (Queue Manager) Clock at 35 MHz.** This clock speed yields a combine incoming and outgoing throughput of 400Mbps for 40cc Enq, Deq commands and 800 Mbps for 20cc Enq, Deq commands.
- **FPGA SRAM Utilization at 95%**
- **FPGA Logic gate Utilization 55%**

Assuming 35MHz as the Queue Manager clock the ABRSU Speed up is 1,4 instead of 2, that was assumed in subsections 7.4.1 - 7.4.3.

In order for the Queue Manager not to be the bottleneck of the system it must be that:

$$\begin{aligned}
 1 \text{ Cell arrival time} &\geq 1 \text{ Enq time} + 1 \text{ Deq time} && \Rightarrow \\
 Tclk_utopia * 28 \text{ utopia_cycles} &\geq Tclk_qm * (20+20) \text{ qm_cycles} && \Rightarrow \\
 Tclk_utopia / Tclk_qm &\geq 50 / 28 = 1,8 && \Rightarrow \\
 40 \text{ ns} / 28\text{ns} &\geq 1,4 && \Rightarrow \\
 1,4 &\geq 1,4 &&
 \end{aligned}$$

This means that the succeeded speed up barely satisfies the Bandwidth needs of the UTOPIA interfaces. If Free-list bypassing and Cell Buffer Pre-allocation were not implemented then the cycles needed for the normal case of enqueue and dequeue would be 20+20 + 10 cycles avoided by Free-list bypassing + 5 cycles avoided by Cell Buffer Pre-allocation = 55 cycles. In that case an ABRSU clock of ~50 MHz of frequency would be necessary.

The combined system performance improvement due to these 2 techniques is 15 cycles / 55 cycles = 27%.

8 Conclusions and Future Work

In this thesis, we studied the architecture of a Per-Flow Queue Manager for the purpose of queueing the ABR traffic of an ATM switch in times of congestion. The Queue Manager was implemented in a large FPGA that was placed on one of the Cards of the Switch. The use of FPGA allowed extensive on board testing of the design during its development and gave us the ability to confirm speed and feasibility assumptions early in the design face.

We used a single SDRAM DIMM memory module for storing of cells and cell pointers. This significantly reduced the pin and trace count of the physical design, yielding a low cost system. The careful scheduling of memory accesses to the SDRAM module by the Queue Manager proved that this single-buffer approach is feasible.

Although dynamic memory allocation increased the number of accesses for each enqueue or dequeue operation and lowered the buffering bandwidth, it allowed us, during testing, to use the Queue Manager for queueing thousands of cells of one flow and maintain the ability to handle the other 64K flows.

The interfacing of the Queue Manager with the external CPU allowed us to both debug the design effectively, insert test traffic that confirmed the correctness of the physical interface of the FPGA with the SDRAM DIMM and the Cubit Pros.

The use of the Free-List bypassing and Cell-Buffer preallocation techniques was proven essential in reaching the goal of near Gbps queueing bandwidth of our design. The 26% improvement of performance they induced with their use compensated for the loss of our clock speed goal (35 MHz instead of the 50 MHz clock that we targeted at the beginning of the design phase). Thus a maximum of 800 Mbps of combined incoming and outgoing ABR throughput was achieved that is sufficient for the ABR traffic needs of a Gbps Switch.

Adding various other switch features to our design is an interesting issue for future work. For instance, the enlargement of the Flow Record from 2 64-bit words to 4 words, could allow the addition of extra fields for each of the 64K supported flows. A New ID field could be added that would be accessible by the external CPU. This field would substitute the Header ID of incoming cells of a Flow. Thus the Queue Manager would offer VP/VC translation along with per-flow queueing.

Another field that could be added, is an Explicit Rate Field. The CPU would compute this appropriate explicit rate and use it to set the Explicit Rate Field inside the RM cells of a given flow, if it is smaller than the RM explicit rate. In that way, the RM explicit rate flow control could be supported, with the software taking care of the rate calculation and the hardware making the RM field enumeration.

Other interesting features that could be supported with extra fields are cell-dropping mechanisms for ill-behaved flows.

9 References

- [1] S. Keshav: "An Engineering Approach to Computer Networking", Addison-Wesley, 1997, ISBN 0-201-63442-2
- [2] R. Lamaire, D. Serpanos: A 2-Dimensional Round-Robin Scheduling Mechanism for Switches with Multiple Input Queues , In IEEE/ACM Transactions on Networking, pages 471-482, 2(5), Oct 1994.
- [3] Y. Oie, M. Murata, K. Kubota and H. Miyahara: «Effect of speedup in non-blocking packet switch,» Proc. ICC '89, Boston, MA, June 1989, pp. 410-414.
- [4] J. G. Jim: «The throughput of data switches with and without speedup», Dai Schools of Industrial and Systems Engineering, and Mathematics Georgia Institute of Technology Balaji Prabhakar.
- [5] A. Charny: «Providing QoS Guarantees in Input Buffered Crossbar Switches with Speedup», PhD Thesis, MIT, 1998.
- [6] P. Prabhakar and N. McKeown: «On the speedup required for combined input-and output-queued switching», to appear in Automatica.
- [7] G. Kornaros, C. Kozyrakis, P. Vatsolaki, M. Katevenis: "Pipelined Multi-Queue Management in a VLSI ATM Switch Chip with Credit-Based Flow Control", in Proc. ARVLSI'97 (17 th Conference on Advanced Research in VLSI), Univ. of Michigan at Ann Arbor, MI USA, Sept. 1997, IEEE Computer Soc. Press, ISBN 0-8186-7913-1, pp. 127-144.
- [8] Ioannis Mavroidis: "Heap Management in Hardware", Technical Report 222, ICS-FORTH, July 1998
- [9] A. Ioannou, M. Katevenis: "Pipelined Heap (Priority Queue) Management for Advanced Scheduling in High Speed Networks", Proc. IEEE Int. Conf. on Communications (ICC'2001), Helsinki, Finland, June 2001, pp. 2043-2047;
<http://archvlsi.ics.forth.gr/mugpro/heapMgt.html>
- [10] A. Nikologiannis, M. Katevenis: "Efficient Per-Flow Queueing in DRAM at OC-192 Line Rate using Out-of-Order Execution Techniques", Proc. IEEE Int. Conf. on Communications (ICC'2001), Helsinki, Finland, June 2001, pp. 2048-2052;
<http://archvlsi.ics.forth.gr/mugpro/queueMgt.html>
- [11] M. Katevenis, D. Serpanos, E. Markatos: "Multi-queue management and scheduling for improved QoS in communication networks", *Proceedings of EMMSEC'97* (European Multimedia Microprocessor Systems and Electronic Commerce Conference), Florence, Italy, Nov. 1997, pp. 906-913;
http://archvlsi.ics.forth.gr/html_papers/EMM-SEC97/paper.html
- [12] Tzi-cker Chiueh, Varadarajan, S.: "Design and evaluation of a DRAM-based shared memory ATM", 1997 ACM International Conference on Measurement and Modeling of Computer Systems (SIGMETRICS 97) Seattle, WA, USA 15-18, June 1997
- [13] P. Andersson, C. Svensson (Lund Univ., Sweden): "A VLSI Architecture for an 80 Gb/s ATM Switch Core", IEEE Innovative Systems in Silicon Conference, Oct. 1996.
- [14] V. Kumar, T. Lakshman, D. Stiliadis: "Beyond Best Effort: Router Architectures for the Differentiated Services of Tomorrow's Internet", IEEE Communications Magazine, May 1998, pp152-164.
- [15] B. Suter, T.V. Lakshman, D. Stiliadis, A.K. Choudhury: "Buffer Management Schemes for Supporting TCP in Gigabit Routers with Per-Flow Queueing", IEEE Journal in Selected Areas in Communications, August 1999.

[16] <http://www.altera.com>

[17] <http://www.motorola.com>

[18] <http://www.transwitch.com>

[19] <http://www.micron.com>

[20] <http://www.pmc-sierra.com>

[21] ATM Forum: "Traffic Management Specification, Version 4.1", AF-TM-0121.000, March 1999.

[22] Ch. Lolas: "Design and Implementation of Low-Level software for high-speed packet switches", Master of Science Thesis, Computer Science Department, University of Crete, Greece, November 2001.

[23] G. Papadakis: "Design and Implementation in FPGA of an ABR Traffic Scheduler for an ATM Switch", Master of Science Thesis, Computer Science Department, University of Crete, Greece, November 2001.