

Εθνικό Μετσόβιο Πολυτεχνείο Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών

Τομέας Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής

Σχεδίαση ενός πλήρως αναλογικού και εξαιρετικά χαμηλής κατανάλωσης βαθέως συνελικτικού νευρωνικού δικτύου



του

Φούφα Ζήση

Επιβλέπων: Παύλος Π. Σωτηριάδης Καθηγητής Ε.Μ.Π.

Αθήνα, Νοέμβριος 2024



Εθνικό Μετσόβιο Πολυτεχνείο Σχολή Ηλεκτρολόγων Μηχανικών και Μηχανικών Υπολογιστών Τομέας Επικοινωνιών, Ηλεκτρονικής και Συστημάτων Πληροφορικής

Σχεδίαση ενός πλήρως αναλογικού και εξαιρετικά χαμηλής κατανάλωσης βαθέως συνελικτικού νευρωνικού δικτύου

Διπλωματική Εργασία

του

Φούφα Ζήση

Επιβλέπων: Παύλος Π. Σωτηριάδης Καθηγητής Ε.Μ.Π.

Εγκρίθηκε από την τριμελή εξεταστική επιτροπή την 29^η Νοεμβρίου 2024:

Παύλος Π. Σωτησιάδης Ναντάριος Κοζύρης Αθαινάσιος Βουλοδάμος

Παύλος Π. Σωτηριάδης Καθηγητής Ε.Μ.Π. Νεκτάριος Κοζύρης Καθηγητής Ε.Μ.Π.

Αθανάσιος Βουλοδήμος Επικ. Καθηγητής Ε.Μ.Π.

Αθήνα, Νοέμβριος 2024

.....

Φούφας Ζήσης

Διπλωματούχος Ηλεκτρολόγος Μηχανικός και Μηχανικός Υπολογιστών, Ε.Μ.Π.

Copyright © Φούφας Ζήσης, 2024. Με επιφύλαξη παντός δικαιώματος. All rights reserved.

Απαγορεύεται η αντιγραφή, αποθήκευση και διανομή της παρούσας εργασίας, εξ ολοκλήρου ή τμήματος αυτής, για εμπορικό σκοπό. Επιτρέπεται η ανατύπωση, αποθήκευση και διανομή για σκοπό μη κερδοσκοπικό, εκπαιδευτικής ή ερευνητικής φύσης, υπό την προϋπόθεση να αναφέρεται η πηγή προέλευσης και να διατηρείται το παρόν μήνυμα. Ερωτήματα που αφορούν τη χρήση της εργασίας για κερδοσκοπικό σκοπό πρέπει να απευθύνονται προς τον συγγραφέα.

Οι απόψεις και τα συμπεράσματα που περιέχονται σε αυτό το έγγραφο εκφράζουν τον συγγραφέα και δεν πρέπει να ερμηνευθεί ότι αντιπροσωπεύουν τις επίσημες θέσεις του Εθνικού Μετσόβιου Πολυτεχνείου.

Περίληψη

Στόχος της παρούσας διπλωματικής είναι η υλοποίηση ενός πλήρως αναλογικού και εξαιρετικά χαμηλής κατανάλωσης βαθέως συνελικτικού νευρωνικού δικτύου. Ειδικότερα αναλύεται η σχεδίαση νέων κυκλωμάτων ηλεκτρονικά ελεγχόμενων αναλογικών πολλαπλασιαστών και κλιμάκωσης εξόδου γενικής χρήσης καθώς και κυκλωμάτων υλοποίησης της συνάρτησης ενεργοποίησης ReLU και του τελεστή argmax. Αχολούθως, συνδυάζοντας τα προαναφερθέντα χυχλώματα χαθώς χαι επίπεδα αναλογικής μνήμης, περιγράφεται η σύνθεση τόσο των συνελικτικών στρωμάτων όσο χαι του πλήρως συνδεδεμένου στρώματος του παραπάνω νευρωνιχού διχτύου ενώ παράλληλα παρέχεται και το πλαίσιο για την αναλογική υλοποίηση παρόμοιων αρχιτεχτονικών. Η διαφορά τάσης τροφοδοσίας του εν λόγω χυχλώματος είναι 0.6V χαι όλα τα ρεύματα που συναντώνται σε αυτό είναι της τάξεως των nA, γεγονός που συνεπάγεται εξαιρετικά χαμηλή κατανάλωση. Το νευρωνικό αυτό δίκτυο εκπαιδεύεται με χρήση της γλώσσας προγραμματισμού Python και το πρωτότυπο μοντέλο σε επίπεδο λογισμικού που προκύπτει χρησιμοποιείται ως σημείο αναφοράς για τις αναλογικές υλοποιήσεις. Εν συνεχεία, το προτεινόμενο αναλογικό μοντέλο αξιολογείται στην ταξινόμηση δειγμάτων, χρησιμοποιώντας ένα κοινό σύνολο δεδομένων μηχανικής όρασης με ποικίλες πρακτικές εφαρμογές. Η σχεδίαση και προσομοίωση όλων των χυχλωμάτων έγινε με χρήση του παχέτου λογισμιχού Cadence IC Suite σε τεχνολογία TSMC 90nm CMOS process.

Λέξεις Κλειδιά: Συνελικτικό Νευρωνικό Δίκτυο, Αναλογικό Νευρωνικό Δίκτυο, On-Chip Classification, Εξαιρετικά Χαμηλή Κατανάλωση, Αναλογικός Πολλαπλασιαστής, Ηλεκτρονικά Ελεγχόμενο Κύκλωμα, Περιοχή Υποκατωφλίου

Abstract

The goal of this diploma thesis is the implementation of a fully analog and ultralow power deep convolutional neural network. More specifically, it presents the design of novel circuits for general purpose electronically controlled analog multiplication and output current scaling as well as circuits for implementing the ReLU activation function and the argmax operator. Subsequently, this thesis describes how through combining the above circuits and layers of analog memory we can synthesize both the convolutional layers and the fully connected layer of the neural network in question while also providing a framework for the analog implementation of similar architectures. The supply voltage difference of this circuit is 0.6V and all currents present in its implementation are in the order of magnitude of a few nA resulting in extremely low power consumption. This neural network was trained using the programming language Python and the resulting prototype software model is used as a reference point for analog implementations. Furthermore, the proposed analog model is evaluated on the correct classification of samples using a common computer vision dataset with many real world applications. The design and simulation of all circuits was conducted using the Cadence IC Suite software package with the TSMC 90nm CMOS process.

Keywords: Convolutional Neural Network, Analog Neural Network, On-Chip Classification, Ultra Low Power, Analog Multiplier, Electronically Controlled Circuit, Subthreshold Region

Ευχαριστίες

Στο σημείο αυτό θα ήθελα να ευχαριστήσω την οικογένεια μου, τους φίλους μου και όλα τα άλλα άτομα που με στήριξαν καθ' όλη την διάρκεια των σπουδών μου.

Θα ήθελα ιδιαίτερα να ευχαριστήσω τον επιβλέποντα καθηγητή μου κ. Παύλο Πέτρο Σωτηριάδη για την πολύτιμη καθοδήγηση του, τις συμβουλές του καθώς και το γνήσιο ενδιαφέρον που επέδειξε κατά την εκπόνηση της διπλωματικής μου εργασίας. Μου δόθηκε η ευκαιρία να ασχοληθώ με ένα πολύ ενδιαφέρον ερευνητικό θέμα σε ένα ιδιαίτερα υποστηρικτικό περιβάλλον. Παράλληλα θα ήθελα να ευχαριστήσω θερμά τον διδάκτορα Βασίλειο Αλιμίση για την ανεκτίμητη βοήθεια του. Η παρούσα εργασία είναι αναμφισβήτητα προϊόν των πολύωρων συζητήσεων και της τακτικής μας επικοινωνίας κατά την διάρκεια των τελευταίων μηνών.

Φούφας Ζήσης, Νοέμβριος 2024

Περιεχόμενα

	Περίληψη	5
	Abstract	6
	Ευχαριστίες	7
	Ευρετήριο Ειχόνων	9
	Κατάλογος Πινάχων	13
1	Εισαγωγή 1.1 Εισαγωγή στη μηχανική μάθηση	14 14 15 17
2	 Επιστημονικό Υπόβαθρο 2.1 Νευρωνικά Δίκτυα 2.1.1 Μοντελο Perceptron 2.1.2 Μοντελο MLP 2.1.3 Βασικές αρχές εκπαίδευσης Νευρωνικών Δικτύων 2.1.4 Συναρτήσεις ενεργοποίησης Νευρωνικών Δικτύων 2.1.5 Συνελικτικά Νευρωνικά Δίκτυα 2.2 MOS τρανζίστορ 2.2.1 Βασική δομή MOS τρανζίστορ 2.2.2 Βασική λειτουργία μεγάλου σήματος MOS τρανζίστορ 2.3 Λειτουργία MOS τρανζίστορ στην περιοχή υποκατωφλίου 	 20 21 23 24 25 26 29 30 32
3	 Δομικά Κυκλώματα 3.1 Κυκλώματα υλοποίησης σιγμοειδούς συνάρτησης 3.1.1 Απλό κύκλωμα σιγμοειδούς με διαφορικό ζεύγος 3.1.2 Κύκλωμα σιγμοειδούς με διαφορικό ζεύγος διαφοράς 3.2 Κυκλώματα πολλαπλασιαστών 	36 36 39 46

	$3.3 \\ 3.4$	3.2.1 Κύκλωμα πολλαπλασιαστή συνελικτικών φίλτρων	46 50 51 52 52 52 54
4	Προ	οτεινόμενη υλοποίηση αναλογικού νευρωνικού δικτύου	62
	4.1	Αρχιτεκτονική προτεινόμενου συνελικτικού νευρωνικού δικτύου	62
	4.2	Κύχλωμα συνελιχτιχού φίλτρου 3×3 ενός χαναλιού εισόδου με ReLU	
		ενεργοποίηση	63
	4.3	Κύκλωμα συνελικτικού φίλτρου $3 imes 3$ τριών καναλιών εισόδου με	
		ReLU ενεργοποίηση	64
	4.4	Υλοποίηση πρώτου επιπέδου με ReLU ενεργοποίηση και average po-	
		oling	66
	4.5	Τλοποίηση επιπέδων 2,3,4 με ReLU ενεργοποίηση	67
	4.6	Τλοποίηση πλήρως συνδεμένου επιπέδου χεφαλής ταξινόμησης	67
5	Πα	ράδειγμα πραγματικής εφαρμογής	72
	5.1	Σύνολο δεδομένων GTSRB	72
	5.2	Εκπαίδευση μοντέλου	72
		5.2.1 Εκπαίδευση μοντέλου σε επίπεδο λογισμικού	73
		5.2.2 Εκπαίδευση μοντέλου σε επίπεδο κυκλωμάτων	73
	5.3	Αποτελέσματα ταξινόμησης συνόλου δεδομένων GTSRB	75
6	Συį	ιπεράσματα και μελλοντική δουλειά	79

Ευρετήριο Ειχόνων

2.1.1 Βασική αρχιτεκτονική μοντέλου Perceptron	21
2.1.2 Παράδειγμα γραμμικά διαχωρίσιμων δεδομένων (αριστερά) και μη γραμ-	
μικά διαχωρίσιμων δεδομένων (δεξιά)	22
2.1.3 Παράδειγμα αρχιτεκτονικής ενός δικτύου MLP με δύο κρυφά επίπεδα	
(οι συνδέσεις των νευρώνων που παραλείπονται θεωρούνται δεδομένες	
ενώ επίσης εννοείται και η ύπαρξη νευρώνων πόλωσης)	23
2.1.4 Κοινές συναρτήσεις ενεργοποίησης νευρωνιχών διχτύων	26
2.1.5 Παράδειγμα δισδιάστατης συνέλιξης δυο πινάχων	28
2.1.6 LeNet 5: Ένα από τα πρώτα CNN που αποτελεί χαραχτηριστικό	
παράδειγμα της τυπικής αρχιτεκτονικής τους	29
2.2.1 Φυσική δομή NMOS τρανζίστορ: Στην επάνω εικόνα παρουσιάζεται	
μια προοπτική άποψη και στην κάτω μια τομή	33
2.2.2 Τομή ολοκληρωμένου CMOS (complementary MOS) κυκλώματος	
όπου για το PMOS τρανζίστορ κατασκευάζεται μια ξεχωριστή περιο-	
χή n ημιαγωγού (n well) στο p υπόστρωμα	34
2.2.3 Μορφή καναλιού NMOS για τρείς περιπτώσεις: (επάνω) $V_{DS} < V_{ov}$,	
(μέση) $V_{DS} = V_{ov}$, (κάτω) $V_{DS} > V_{ov}$	34
2.2.4 Βασικές περιοχές λειτουργίας NMOS τρανζίστορ	35
2.2.5 Εξάρτηση I_D NMOS τρανζίστορ από την V_{GS} στις περιοχές υποκα-	
τωφλίου και κορεσμού	35
3.1.1 Σνεδιάνοαμμα απλού χυχλώματος συγμοειδούς συνάστασας	37
3.1.2 Παράδεινας σξόδου σπλού χυχλόματος σιγμοείδούς για Vr = 0. Ibias	01
-10n A	38
$313 \Pi_{\alpha\alpha}$	00
5.1.5 Παραθεί μα εξόδου απού χολοφατός στημοείδους συναρτήσει του Ibias νια Vr = 0	39
314Σ	05
τιχές τιμές Ibias	40
315 Σχεδιάχοαμμα χυχλώματος σινμοειδούς συνάστησης με διαφορικό $\zeta_{\text{F-}}$	10
ύνος διαφοράς	41
······································	

3.1.6 Εξάρτηση εξόδου ρεύματος του κυκλώματος της εικόνας 3.1.5 από την τάση Vc	42
3.1.7 Παράδειγμα εξόδου χυχλώματος του Σχήματος 3.1.5 συναρτήσει του Ibias για Vr = 0	46
3.1.8 Συνάρτηση μεταφοράς χυχλώματος του Σχήματος 3.1.5 για διαφορετικές τιμές Ibias	47
3.1.9 Κανονικοποιημένες τιμές των $I_{D12}, I_{D14}, Iout$ για διαφορετικές τιμές V_c όπως προκύπτουν απο τις αναλυτικές σχέσεις (3.1.4) και (3.1.5).	48
3.1.10 χεδιάγραμμα χυχλώματος προσήμου	49
μα προσήμου	50
3.2.1 Σχεδιάγραμμα κυκλώματος πολλαπλασιαστή συνελικτικών φίλτρων (Κάθε ένα από τα μπλέ ορθογώνια χρησιμοποιείται στη θέση του κυκλώματος του Σχήματος 3.1.11 που βρίσκεται εντός του μωβ περι-	
γράμματος)	55
διαγραμμα εστιαζει στις χαμηλες τιμες III για να φανει η σημαντικα βελτιωμένη γραμμικότητα αυτού του κυκλώματος)	56
διαφορετικές τιμές Ιin, για θετικό (επάνω) και αρνητικό (κάτω) πρόσημο	57
3.2.4 Σχεδιάγραμμα κυκλώματος κλιμάκωσης εξόδου ρεύματος	58
3.3.1 Σχεδιάγραμμα χυχλώματος ReLU συνάρτησης ενεργοποίησης	59
3.4.1 Σχεδιάγραμμα Lazzaro WTA χυχλώματος Ν εισόδων	59
3.4.2 Σχεδιαγραμμα Lazzaro WTA κυκλωματος 2 εισόδων	60
για $Iin1 = 5nA$ και $Ibias_wta = 10nA$	60
3.4.4 Σχεδιάγραμμα κασκοδικού WTA κυκλώματος Ν εισόδων	61
4.1.1 Αρχιτεκτονική προτεινόμενου συνελικτικού νευρωνικού δικτύου	64
4.2.1 Σχεοιαγραμμα χυκλωματός 5 × 5 συνελιχτιχού φιλτρού ενός καναλίου με ReLU ενεργοποίηση	65
4.2.2 Παράδειγμα υπολογισμών κυκλώματος 3 × 3 συνελικτικού φίλτρου ενός καναλιού	66
4.3.1 Σχεδιάγραμμα κυκλώματος 3 × 3 συνελικτικού φίλτρου τριών κανα- λιών με ReLU ενεργοποίηση	69
4.3.2 Παράδειγμα υπολογισμών κυκλώματος 3 × 3 συνελικτικού φίλτρου τριών καναλιού	70
4.4.1 Σχεδιάγραμμα χυχλώματος πρώτου επιπέδου με ReLU ενεργοποίηση και average pooling	70
4.6.1 Σχεδιάγραμμα κυκλώματος πλήρως συνδεδεμένου επιπέδου κεφαλής	71

5.2.1 Έξοδος ρεύματος επιπέδου average pooling για το δεύτερο και το	
τρίτο φίλτρο του πρώτου συνελικτικού επιπέδου χωρίς προσαρμογή	
των Ibias, Vsm, Vdm	77
5.2.2 Έξοδος ρεύματος επιπέδου average pooling για το δεύτερο και το	
τρίτο φίλτρο του πρώτου συνελιχτιχού επιπέδου μετά απο προσαρμογή	
των Ibias, Vsm, Vdm	78

Κατάλογος Πινάχων

3.1	Διαστάσεις τρανζίστορ (Σχήμα 3.1.5)	44
3.2	Διαστάσεις τρανζίστορ (Σχήμα 3.1.10)	45
3.3	Διαστάσεις επιπρόσθετων τρανζίστορ Σχήματος 3.2.1	47
3.4	Διαστάσεις τρανζίστορ Σχήματος 3.3.1	51
5.1	Αχρίβεια ταξινόμησης 1000 δειγμάτων του συνόλου δεδομένων GTSRB	76
5.2	Μέγιστη κατανάλωση ενέργειας μετρημένη για 250 δείγματα του συ-	
	νόλου δεδομένων GTSRB	76

Κεφάλαιο 1

Εισαγωγή

1.1 Εισαγωγή στη μηχανική μάθηση

Ο όρος μηχανική μάθηση αναφέρεται στο πεδίο εκείνο της τεχνητής νοημοσύνης που μελετάει τη χρήση δεδομένων και αλγορίθμων για την επίτευξη διαφόρων εργασιών από μηχανές με τρόπο παρόμοιο με έναν άνθρωπο [2]. Ειδικότερα οι τεχνικές μηχανικές μάθησης επιτρέπουν τη δημιουργία υπολογιστικών μοντέλων που δεδομένου ενός συνόλου δεδομένων και κανόνων και ενός αλγορίθμου ανανέωσης των παραμέτρων τους, έχουν τη δυνατότητα να ενσωματώνουν προηγούμενη γνώση για την βελτίωση της επίδοσης τους σε μια εργασία, χωρίς περεταίρω ανθρώπινη παρέμβαση [3]. Έτσι τα μοντέλα αυτά μπορούν να μάθουν να λύνουν ένα πρόβλημα χωρίς να προγραμματιστούν από την αρχή με τα ακριβή βήματα για την επίλυση του. Πρόσφατα, ο τομέας της μηχανικής μάθησης έχει παρουσιάσει σημαντική άνθιση λόγω της ανάπτυξης νέων αλγορίθμων και αρχιτεκτονικών μοντέλων, αλλά και, σε μεγάλο βαθμό, λόγω του τεράστιου όγκου δεδομένων που έχει παραχθεί τα τελευταία χρόνια, καθώς και της σημαντικής μείωσης του κόστους των υπολογιστικών πόρων [7]

Οι υλοποιήσεις μηχανικής μάθησης χωρίζονται σε τρείς βασικές κατηγορίες: επιβλεπομένη μάθηση, μη-επιβλεπόμενη μάθηση και ενισχυτική μάθηση [7]. Η επιβλεπόμενη μάθηση, που είναι και ο πιο συχνά χρησιμοποιούμενος τύπος μηχανικής μάθησης, αναφέρεται στην εκπαίδευση ενός μοντέλου με τη χρήση ενός συνόλου δεδομένων, όπου για κάθε διάνυσμα εισόδου εκπαίδευσης είναι γνωστή η επιθυμητή τιμή εξόδου του μοντέλου. Η επιβλεπόμενη μάθηση χρησιμοποιείται συχνά για προβλήματα όπως η αναγνώριση εικόνων και η πρόβλεψη χρονοσειρών. Η μη επιβλεπόμενη μάθηση περιγράφει την διαδικασία εκείνη μάθησης κατά την οποία δεν ορίζεται συγκεκριμένη επιθυμητή έξοδος του μοντέλου και απλώς παρέχονται σε αυτό κάποιες υποθέσεις για το σύνολο δεδομένων [7]. Έτσι ένα μοντέλο μηχανικής μάθησης μπορεί να χρησιμοποιηθεί για να ανακαλύψει άγνωστα μοτίβα που υπάρχουν στα δεδομένα εκπαίδευσης. Παραδείγματα τέτοιων προβλημάτων είναι η κατάτμηση εικόνων σε ομοιόμορφες περιοχές (image segmentation) και η μείωση διαστάσεων (dimensionality reduction), όπου στόχος είναι η εύρεση μιας αναπαράστασης ενός συνόλου δεδομένων με λιγότερες διαστάσεις. Η ενισχυτική μάθηση αναφέρεται στην διαδικασία κατά την οποία ένα μοντέλο επιχειρεί να ανακαλύψει τη βέλτιστη στρατηγική για ένα πρόβλημα, αλληλεπιδρώντας με το «περιβάλλον» του και προσπαθώντας να μεγιστοποιήσει ένα προκαθορισμένο κριτήριο. Στην περίπτωση αυτή η πληροφορία που παρέχεται από το σύνολο δεδομένων είναι μεταξύ αυτής των περιπτώσεων της επιβλεπόμενης και της μη-επιβλεπόμενης μάθησης, καθώς για κάθε έξοδο του μοντέλου παρέχεται μόνο μία ένδειξη ως προς το αν αυτή είναι η επιθυμητή. Η ενισχυτική μάθηση χρησιμοποιείται για εφαρμογές όπως η ρομποτική και η αυτόνομη οδήγηση [8]. Τέλος, αντικείμενο έντονης μελέτης είναι και ο συνδυασμός των παραπάνω τύπων μηχανικής μάθησης, όπως στην περίπτωση της ημι-επιβλεπόμενης μάθησης, όπου χρησιμοποιείται ένα μικρό σύνολο σεσημασμένων δεδομένων (labeled data) για να κατευθύνει την διαδικασία εκπαίδευσης σε ένα μεγάλο σύνολο μη-σεσημασμένων δεδομένων.

1.2 Εισαγωγή στα συνελικτικά νευρωνικά δίκτυα

Τα τεχνητά νευρωνικά δίκτυα είναι μια κατηγορία μοντέλων επιβλεπόμενης μηχανικής μάθησης των οποίων η λειτουργία εμπνέεται από αυτήν των βιολογικών νευρώνων. Το ανθρώπινο νευρικό σύστημα περιλαμβάνει κύτταρα που ονομάζονται νευρώνες χαι συνδέονται μεταξύ τους με την χρήση νευραξόνων χαι δενδριτών [1]. Ένας βιολογιχός νευρώνας δέχεται εισόδους από τους δενδρίτες του χαι μεταφέρει το σήμα εξόδου του μέσω του νευράξονα που διαθέτει. Συγκεκριμένα, αυτός δέχεται το συνολικό σήμα εισόδου από τις περιοχές των δενδριτών, μέσω της έκλυσης χημικών μορίων που ονομάζονται νευροδιαβιβαστές. Εφόσον το ηλεκτρικό δυναμικό της κυτταριχής του μεμβράνης ξεπεράσει μια τιμή χατωφλίου, πυροδοτείται η έχλυση ενός σήματος δυναμικού ενέργειας το οποίο ταξιδεύει κατά μήκος του νευράξονα. Αντιθέτως, σε περίπτωση που το δυναμικό της κυτταρικής μεμβράνης είναι μικρότερο της τιμής χατωφλίου, ο νευρώνας παραμένει σε χατάσταση ηρεμίας [9]. Στη συνέχεια οι απολήξεις των νευραξώνων πολλών νευρώνων συνδέονται με δενδρίτες άλλων νευρώνων σε περιοχές που ονομάζονται συνάψεις. Οι ισχύεις αυτών των συναπτικών συνδέσεων είναι μεταβλητές και μπορούν να αλλάξουν με βάση εξωτερικά ερεθίσματα, οδηγώντας έτσι στην εχμάθηση από τον άνθρωπο νέων πληροφοριών με βάση την εμπειρία. Με ανάλογο τρόπο ορίζονται και οι βασικές ιδέες πίσω από την αρχιτεχτονιχή των τεχνητών νευρωνιχών διχτύων. Έτσι τα τεχνητά νευρωνιχά δίχτυα περιέχουν μονάδες υπολογισμού που ονομάζονται επίσης νευρώνες και συνδέονται μεταξύ τους με αχμές μεταβλητών βαρών που λειτουργούν χατά αναλογία των συναπτιχών συνδέσεων των βιολογιχών νευρώνων. Αντίστοιχα η έξοδος χάθε τεχνητού νευρώνα καθορίζεται από το άθροισμα των σημάτων που καταλήγουν σε αυτόν και μια προκαθορισμένη συνάρτηση ενεργοποίησης η οποία μπορεί να ποικίλει ανάλογα την συμπεριφορά του δικτύου που θέλουμε να επιτύχουμε. Η εκπαίδευση ενός τέτοιου μοντέλου γίνεται κατ'αντιστοιχία με τους βιολογικούς νευρώνες μέσω του ερεθίσματος που παρέχεται από το σύνολο δεδομένων εκπαίδευσης το οποίο περιέχει παραδείγματα επιθυμητών ζευγών εισόδων-εξόδων του δικτύου [1]. Πιο συγκεκριμένα σε κάθε βήμα εκπαίδευσης, η έξοδος του δικτύου συγκρίνεται με την επιθυμητή έξοδο για τις εκάστοτε εισόδους και υπολογίζεται το σφάλμα με βάση μια προκαθορισμένη συνάρτηση σφάλματος. Ακολούθως, αυτό τροφοδοτείται, μέσω κατάλληλου αλγορίθμου, στο νευρωνικό δίκτυο, ώστε με βάση αυτήν την νέα πληροφορία να ανανεωθούν κατάλληλα τα βάρη των συνδέσεων μεταξύ των νευρώνων.

Σήμερα τα τεχνητά νευρωνικά δίκτυα, τα οποία στην παρούσα διπλωματική εργασία θα αποκαλούνται στο εξής απλώς νευρωνικά δίκτυα χάριν συντομίας, έχουν επικρατήσει ως ένας από τους πιο ευρέως χρησιμοποιούμενους τύπους μοντέλων μηχανικής μάθησης. Το γεγονός αυτό συνδέεται άμεσα με την μεγάλη ευελιξία τους και την ικανότητα τους να αφομοιώνουν αποδοτικά, με τις κατάλληλες αρχιτεκτονικές, γνώση από πολύ μεγάλα σύνολα δεδομένων. Έτσι, μοντέλα νευρωνικών δικτύων χρησιμοποιούνται σε πληθώρα σύγχρονων εφαρμογών, όπως η πρόβλεψη χρονοσειρών, η όραση υπολογιστών η επεξεργασία φυσικής γλώσσας και η ανάλυση βιολογικών σημάτων.

Στο πλαίσιο της πολύ μεγάλης ανάπτυξης του επιστημονιχού χλάδου των νευρωνιχών διχτύων τα τελευταία χρόνια, ερευνήθηχε η χρήση ποιχίλων σχετιχών αρχιτεχτονιχών. Μια από τις πιο επιτυχημένες τέτοιες αρχιτεχτονιχές είναι αυτή των συνελικτικών νευρωνικών δικτύων η οποία εμπνεύστηκε σε μεγάλο βαθμό από τα πειράματα των Hubel και Weisel για την οργάνωση των νευρώνων στον οπτικό φλοιό της γάτας [13]. Πιο συγκεκριμένα τα συνελικτικά νευρωνικά δίκτυα (convolutional neural networks ή CNNs) είναι ένας τύπος νευρωνιχών διχτύων χωρίς ανάδραση (feedforward neural networks) που χρησιμοποιούν συνελικτικά φίλτρα αντί για απλούς νευρώνες σε ένα τουλάχιστον στρώμα τους. Τα νευρωνικά δίκτυα αυτού του τύπου είναι σχεδιασμένα να επεξεργάζονται δεδομένα με γνωστή πλεγματική δομή [6] και ενσωματώνουν στην αρχιτεκτονική τους την υπόθεση ότι τα δεδομένα εισόδου τους παρουσιάζουν χωρική συσχέτιση [11]. Έτσι τα συνελικτικά νευρωνικά δίκτυα έχουν εφαρμοστεί με μεγάλη επιτυχία σε προβλήματα αναγνώρισης εικόνων [11, 14, 15], αλλά η χρησιμότητα τους επεχτείνεται σε προβλήματα με χάθε είδους δεδομένα με χωρική συσχέτιση, όπως κείμενα και χρονικές σειρές [1]. Ακόμη τα CNNs έχουν γενικά πολύ μικρότερο αριθμό παραμέτρων και συνδέσεων από ένα συγκρίσιμο σε μέγεθος απλό νευρωνικό δίκτυο χωρίς ανάδραση (όπως επεξηγείται στην ενότητα 2.1.5). Αυτό συνεπάγεται πως η εκπαίδευση τους είναι υπολογιστικά ευκολότερη, χωρίς τα δύο να έχουν σημαντικές διαφορές στην θεωρικά βέλτιστη επίδοση τους [11].

Χρησιμότητα αναλογικής υλοποίησης μοντέλων Μηχανικής Μάθησης

Τα μοντέλα μηχανικής μάθησης έχουν σήμερα κρίσιμο ρόλο στην εξαγωγή χρήσιμης πληροφορίας από τον τεράστιο όγχο δεδομένων που παράγεται χάθε μέρα από εχατομμύρια διαφορετιχούς αισθητήρες ανά τον χόσμο. Μεγάλο μέρος όμως της επεξεργασίας αυτών των δεδομένων γίνεται στο νέφος (cloud), σε servers πολύ μαχριά από τους ίδιους τους αισθητήρες, γεγονός που απαιτεί την μεταφορά μεγάλου όγχου δεδομένων μέσω του διχτύου. Λόγω αυτού του μεγάλου όγχου δεδομένων αλλά και πιθανών περιορισμών στο διαθέσιμο εύρος ζώνης του δικτύου, αυτό το μοντέλο χαθίσταται ιδιαίτερα αχριβό χαι ενεργοβόρο, ενώ παράλληλα παρουσιάζει επιπλέον μειονεκτήματα ασφάλειας και ταχύτητας [5]. Μια λύση σε πολλά από τα προβλήματα που παρουσιάζει η επεξεργασία «raw» δεδομένων απευθείας στο νέφος δίνει η ιδέα της χρήσης ενσωματωμένων υπολογιστικών συστημάτων κοντά στον αισθητήρα. Έτσι ένα μεγάλο μέρος της απαιτούμενης επεξεργασίας μπορεί να λαμβάνει χώρα τοπικά και στην συνέχεια, εφόσον αυτό είναι αναγκαίο, να μεταφέρεται στο νέφος ένας μικρός, συγκριτικά με τον αρχικό, όγκος δεδομένων για την διενέργεια πιο σύνθετων υπολογισμών η την αποθήχευση του για μετέπειτα χρήση. Αχόμη, ενσωματωμένα συστήματα που μπορούν να υλοποιούν σχετικά απλά μοντέλα μηχανικής μάθησης, αλλά και διαφόρου είδους υπολογισμούς γενικού σκοπού είναι ιδιαίτερα χρήσιμα σε εφαρμογές όπως εμφυτευμένες ιατρικές συσκευές [16, 17, 18] και αισθητήρες σε πολύ απομαχρυσμένες περιοχές που χρειάζεται να λειτουργούν χωρίς σύνδεση στο δίκτυο και χωρίς κάποια ισχυρή πηγή ενέργειας. Έτσι καθίσταται προφανές πως οι ενσωματωμένες συσκευές που περιγράφηκαν προηγουμένως χρειάζεται να χαρακτηρίζονται από μεγάλη διεκπεραιωτικότητα, χαμηλό ποσοστό σφαλμάτων, χαμηλό κόστος και ιδιαίτερα χαμηλή κατανάλωση ενέργειας. Παράλληλα, συχνά απαιτείται χαι η δυνατότητα προσαρμογής των παραμέτρων τους, ώστε αυτές να μπορούν να χρησιμοποιηθούν για πολλές διαφορετικές εφαρμογές και υπό διαφορετικές συνθήκες [5].

Σήμερα, η πλειοψηφία των ενσωματωμένων συστημάτων υλοποιούνται με την χρήση ψηφιαχών χυχλωμάτων. Τα συστήματα αυτά μπορούν να βασίζουν την λειτουργία τους σε κάποιον τύπο επεξεργαστή γενιχής χρήσης, όπως CPUs και GPUs ή σε πιο εξειδιχευμένα χυχλώματα όπως FPGAs ή ψηφιαχά ASIC. Τα FPGAs (field programmable gate arrays) είναι ψηφιαχά χυχλώματα αποτελούμενα από προγραμματιζόμενα μπλοχ λογικής, τα οποία μπορούν να προγραμματιστούν πολλαπλές φορές από τον χρήστη για να επιτελούν μια συγχεχριμένη λειτουργία [21]. Τα ASICs (application specific integrated circuits) είναι χυχλώματα τα οποία σχεδιάζονται και κατασχευάζονται εξ' αρχής για να επιτελούν μια συγχεχριμένη λειτουργία. Το γεγονός αυτό τα καθιστά ιδιαίτερα γρήγορα και αποδοτικά, αλλά και αρχετά αχριβά στον σχεδιασμό. Τέτοια εξειδιχευμένα ψηφιαχά χυχλώματα έχουν χρησιμοποιηθεί με αρχετή επιτυχία ως επιταχυντές υλιχού (hardware accelerators) για εφαρμογές μηχανικής μάθησης [19, 20, 23], προσφέροντας σημαντική βελτίωση απόδοσης έναντι της χρήσης επεξεργαστικών μονάδων γενικής χρήσης. Όμως, ενώ τα προαναφερθέντα εξειδιχευμένα επεξεργαστιχά συστήματα για εφαρμογές μηχανιχής μάθησης υλοποιούνται κατά κύριο λόγο με την χρήση ψηφιακών κυκλωμάτων, ακόμη μεγαλύτερη απόδοση μπορεί να επιτευχθεί με την χρήση αναλογικών η μικτού σήματος αρχιτεχτονιχών [25]. Μια τέτοια προσέγγιση είναι ιδιαίτερα ελχυστιχή για ML (Machine Learning) εφαρμογές, καθώς αυτές συγνά παρουσιάζουν υψηλή ανογή στον θόρυβο και την αβεβαιότητα [24], σε αντίθεση με εφαρμογές όπως οι ηλεκτρονικές τραπεζικές συναλλαγές ή πολλές ιατρικές εξετάσεις. Αυτό συνάδει και με το γεγονός ότι στην περίπτωση των νευρωνιχών διχτύων οι δύο αυτές έννοιες είναι βασιχά χαρακτηριστικά του τρόπου εκπαίδευσης τους. Έτσι αυτά τα μοντέλα μπορούν να ωφεληθούν από το εγγενές πλεονέχτημα των αναλογιχών χυχλωμάτων, αυτό δηλαδή της αυξημένης ταχύτητας και μειωμένης κατανάλωσης με κόστος την υπολογιστική αχρίβεια. Παράλληλα, λόγω της πολύ χαμηλής χατανάλωσης τους, αναλογικά χυκλώματα που εκτελούν inference (εξαγωγή συμπερασμάτων για ένα σύνολο εισόδων) μοντέλων μηχανιχής μάθησης μπορούν να χρησιμοποιηθούν ως wake-up circuits. Έτσι μπορούν να μειώσουν σημαντικά την κατανάλωση ενέργειας ενός συστήματος που περιλαμβάνει πιο σύνθετα ψηφιακά κυκλώματα τα οποία είναι χρήσιμα για την υψηλή τους αχρίβεια σε χάποιες περιπτώσεις [16]. Ως wake-up circuit ορίζεται ένα ολοχληρωμένο χύχλωμα - αναλογικό σε αυτήν την περίπτωση - το οποίο θα είναι μονίμως ανοιχτό, χωρίς να επηρεάζει την συνολιχή χατανάλωση ενός συστήματος. Ο σκοπός του είναι να ενεργοποιεί ένα πιο σύνθετο και ενεργοβόρο κύκλωμα σε περίπτωση που οι συνθήχες το απαιτήσουν, όπως στην περίπτωση του [16] όπου ένα αναλογικό κύκλωμα προσπαθεί συνεχώς να προβλέψει αν ένας ασθενής θα πάθει επιληπτιχό επεισόδιο χαι σε περίπτωση πιθανής πρόβλεψης ενεργοποιεί ένα ψηφιαχό κύκλωμα για να κρίνει την ορθότητας αυτής.

Τα τελευταία χρόνια έχει ερευνηθεί η αναλογική υλοποίηση αρκετών διαφορετικών μοντέλων μηχανικής μάθησης. Αυτά συμπεριλαμβάνουν κλασικά συστήματα με λίγες παραμέτρους όπως Bayesian ταξινομητές [16, 17] και υλοποιήσεις του αλγορίθμου Support Vector Machine [29] αλλά και υλοποιήσεις διαφόρων μοντέλων νευρωνικών δικτύων. Η πλειοψηφία των αναλογικών και υβριδικών επιταχυντών για νευρωνικά δίκτυα κάνουν χρήση υπολογισμών κοντά/εντός μνήμης (in memory/near memory computing) σε συνδυασμό με μείωση της υπολογιστικής ακρίβειας [29]. Τέτοια παραδείγματα περιλαμβάνουν τις υλοποιήσεις [31, 32, 33] που χρησιμοποιούν κελιά SRAM και τις υλοποιήσεις [34, 35] που χρησιμοποιούν τις αναλογικές ιδιότητες μη πτητικής μνήμης για να εκτελέσουν ταυτόχρονα πολλούς υπολογισμούς MAC (multiply and accumulate). Τα παραπάνω συστήματα στοχεύουν κατά κύριο λόγο στη βελτίωση της ταχύτητας και της κατανάλωσης ενέργειας σε γενικής χρήσεως επεξεργαστικές μονάδες που χρησιμοποιούνται κατά την εκπαίδευση και την χρήση νευρωνικών δικτύων. Άλλες προσεγγίσεις κάνουν χρήση ξεχωριστών κυκλωμάτων πολλαπλασιαστών, όπως το [36] όπου χρησιμοποιούνται ζευγάρια binary-weighted current steering digital-to-analog converters (DACs) για να υλοποιήσουν μία επεξεργαστική μονάδα συνελικτικών δικτύων με 27 MAC κυκλώματα και κατανάλωση ισχύος 540.6µW. Ακόμη δύο παραδείγματα αποτελούν το [37] που χρησιμοποιεί ένα Gilbert-cell πολλαπλασιαστή και έναν ενισχυτή μεταβλητού κέρδους καθώς και το [38] το οποίο χρησιμοποιεί αναλογικούς πολλαπλασιαστές που βασίζονται στην translinear αρχή.

Κεφάλαιο 2

Επιστημονικό Υπόβαθρο

Στο παρόν χεφάλαιο αναλύεται το απαιτούμενο επιστημονιχό υπόβαθρο για την κατανόηση της εν λόγω διπλωματιχής εργασίας. Στο πρώτο μέρος γίνεται μια εχτενής παρουσίαση των βασιχών αλγορίθμων χαι εννοιών γύρω από το πεδίο των νευρωνιχών διχτύων. Αυτή ξεχινάει από τις απλούστερες δυνατές αρχιτεχτονιχές και χτίζοντας επάνω στις ιδέες που αυτές εισάγουν οδηγούμαστε σταδιαχά στα βαθιά συνελιχτιχά νευρωνιχά δίχτυα, όπως το μοντέλο που υλοποιείται στην παρούσα εργασία. Στο δεύτερο μέρος αναλύεται σε βάθος η λειτουργία του βασιχού χυχλωματιχού δομιχού στοιχείου του εν λόγω αναλογιχού χυχλώματος, όπως χαι της πλειοψηφίας των σύγχρονων ολοχληρωμένων χυχλωμάτων, το MOS τρανζίστορ. Συγχεχριμένα, παρουσιάζονται οι βασιχές αρχές της δομής χαι της λειτουργίας μεγάλου σήματος χαθώς και αυτές της λειτουργίας στην περιοχή υποχατωφλίου η οποία προσφέρει σημαντιχά πλεονεχτήματα για την συγχεχριμένη εφαρμογή.

2.1 Νευρωνικά Δ ίκτυα

Όπως αναφέρθηκε και προηγουμένως, τα νευρωνικά δίκτυα είναι ένας τύπος μοντέλων επιβλεπόμενης μηχανικής μάθησης που εμπνέεται από την λειτουργία του ανθρώπινου εγκεφάλου και λειτουργεί χρησιμοποιώντας υπολογιστικές μονάδες (νευρώνες) που συνδέονται μεταξύ τους μέσω ακμών μεταβλητού βάρους. Σε αυτήν την ενότητα παρουσιάζεται αρχικά ο απλούστερος τύπος ενός νευρωνικού δικτύου, το μοντέλο Perceptron, που διαθέτει ένα μόνο επίπεδο, ακολουθούμενο από το μοντέλο MLP (Multi-Layer Perceptron). Στη συνέχεια, αναλύονται οι βασικές ιδέες γύρω από τις τεχνικές εκπαίδευσης των νευρωνικών δικτύων ενώ παράλληλα γίνεται μια ανασκόπηση των διαφορετικών συναρτήσεων ενεργοποίησης και της επίδρασης που αυτές μπορούν να έχουν. Ακολούθως, έχοντας χτίσει το απαραίτητο υπόβαθρο, αναλύονται οι κεντρικές ιδέες των συνελικτικών νευρωνικών δικτύων.

2.1.1 Μοντελο Perceptron

Το μοντέλο Perceptron είναι η απλούστερη δυνατή μορφή ενός νευρωνικού δικτύου. Αυτό περιλαμβάνει ένα επίπεδο εισόδων περιλαμβάνει d κόμβους οι οποίοι δέχονται ως εισόδους τα d χαρακτηριστικά μιας δεδομένης εισόδου $\bar{X} = [x_1, x_2, ..., x_d]$ και τα μεταδίδουν μέσω ακμών με βάρη $\bar{W} = [w_1, w_2, ..., w_d]$ στον κόμβο εξόδου. Ακόμη το μοντέλο αυτό περιλαμβάνει και έναν νευρώνα πόλωσης ο οποίος μεταδίδει πάντα μια σταθερή τιμή. Ακολούθως ο πολλαπλασιασμός των χαρακτηριστικών εισόδου με τα αντίστοιχα βάρη και η άθροιση τους $\bar{W} \cdot \bar{X} = \sum_{i=1}^{d} w_i x_i$ λαμβάνει χώρα στον κόμβο εξόδου, όπου σε αυτό το αποτέλεσμα αθροίζεται και η συνεισφορά του νευρώνα πόλωσης. Στην συνέχεια, η τελική πρόβλεψη εξάγεται με βάση το πρόσημο αυτής της ποσότητας ως εξής:

$$\hat{y} = sign\{\bar{W} \cdot \bar{X} + b\} = sign\{\sum_{i=1}^{d} w_i x_i + b\}$$

όπου η περισπωμένη πάνω από το y υποδειχνύει ότι πρόχειται για πρόβλεψη του μοντέλου χαι όχι παρατηρούμενη τιμή. Δεδομένου τώρα ότι η συνάρτηση πρόσημου, που σε αυτήν την περίπτωση χρησιμοποιείται ως συνάρτηση ενεργοποίησης του νευρώνα εξόδου, μπορεί να πάρει μόνο δύο διαχριτές τιμές {-1, 1}, οι εφαρμογές αυτού του συστήματος περιορίζονται σε προβλήματα δυαδιχής ταξινόμησης.





Το σφάλμα αυτής της πρόβλεψης, που ορίζεται ως $E(X) = y - \hat{y}$, λαμβάνει τιμές στο σύνολο -2, 0, 2. Για να μπορεί το συγχεχριμένο δίχτυο να βελτιώνει την αχρίβεια πρόβλεψης του χρησιμοποιώντας τα δεδομένα του συνόλου εχπαίδευσης, θα πρέπει σε περίπτωση λανθασμένης πρόβλεψης, η χατεύθυνση ενημέρωσης των βαρών των αχμών του να είναι αντίθετη αυτής του σφάλματος. Συγχεχριμένα όταν η είσοδος $ar{X}$ εισάγεται στο δίκτυο, το διάνυσμα βαρών $ar{W}$ θα ενημερώνεται με βάση τον εξής αλγόριθμο κατάβασης δυναμικού:

$$\bar{W} \leftarrow \bar{W} + aE(\bar{X})\bar{X}$$

(perceptron learning rule/delta rule) όπου το α είναι μια παράμετρος που ελέγχει τον ρυθμό μάθησης του μοντέλου.

Όταν ο αλγόριθμος Perceptron παρουσιάστηκε για πρώτη φορά από τον Rosenblatt [39], αυτός στόχευε στην ελαχιστοποίηση του σφάλματος χωρίς να έχει παρουσιαστεί μια τυπική διατύπωση της βελτιστοποίησης αυτής. Όμως αυτό το πρόβλημα ελαχιστοποίησης μπορεί να διατυπωθεί σε μορφή ελαχίστων τετραγώνων σε ένα σύνολο δεδομένων εκπαίδευσης D ως εξής:

$$Minimize_{\bar{W}}L = \sum_{(\bar{X},y)\in D} (y-\hat{y})^2 = \sum_{(\bar{X},y)\in D} (y-sign\{\bar{W}\cdot\bar{X} + b\})^2$$

Ο τύπος του μοντέλου που προτείνεται στο Perceptron είναι ένα γραμμικό μοντέλο στο οποίο η εξίσωση $\overline{W} \cdot \overline{X} + b = 0$ ορίζει ένα γραμμικό υπερεπίπεδο. Ένα μοντέλο τέτοιου τύπου συγκλίνει πάντα σε μία λύση με μηδενικό σφάλμα όταν τα δεδομένα εκπαίδευσης είναι γραμμικά διαχωρίσιμα [39] (όπως φαίνεται στο παράδειγμα του σχήματος 2.1.2) αλλά όχι σε αντίθετη περίπτωση κατά την οποία μπορεί να καταλήξει σε μία πολύ κακή λύση.

Σχήμα 2.1.2: Παράδειγμα γραμμικά διαχωρίσιμων δεδομένων (αριστερά) και μη γραμμικά διαχωρίσιμων δεδομένων (δεξιά)



Γενικά, παρά το γεγονός ότι το Perceptron προσομοιάζει αρκετά παραδοσιακά μοντέλα μηχανικής μάθησης, η θεώρηση του ως υπολογιστική μονάδα επιτρέπει την δημιουργία πολύ ισχυρότερων συστημάτων συγκεντρώνοντας μαζί πολλούς τέτοιους νευρώνες.

2.1.2 Μοντελο MLP

Όπως αναφέρθηκε προηγουμένως, το μοντέλο Perceptron αδυνατεί να παράγει καλές λύσεις όταν τα δεδομένα του προβλήματος δεν είναι γραμμικά διαχωρίσιμα, γεγονός που ισχύει ακόμη και όταν χρησιμοποιούνται άλλες συναρτήσεις ενεργοποίησης πέραν της συνάρτησης προσήμου. Έτσι για την αντιμετώπιση τέτοιων προβλημάτων προτάθηκε η χρήση μίας πιο σύνθετης αρχιτεκτονικής νευρωνικών δικτύων, του Multi-Layer Perceptron. Η συγκεκριμένη αρχιτεκτονική περιλαμβάνει ένα η περισσότερα υπολογιστικά επίπεδα μεταξύ των επιπέδων εισόδου και εξόδου, τα οποία χαρακτηρίζονται ως κρυφά επειδή οι εκτελούμενοι σε αυτά υπολογισμοί δεν είναι άμεσα ορατοί στον χρήστη. Αυτά τα μοντέλα αναφέρονται συχνά και ως δίκτυα εμπροσθόδοσης δεδομένων (feedforward networks) επειδή τα επίπεδα τους τροφοδοτούνται διαδοχικά από το προηγούμενο στο επόμενο, με κατεύθυνση απο την είσοδο προς το επίπεδο εξόδου, ενώ σε αυτά κάθε κόμβος ενός επιπέδου συνδέεται με όλους τους κόμβους του επόμενου (πλήρως συνδεδεμένα επίπεδα).

Σχήμα 2.1.3: Παράδειγμα αρχιτεκτονικής ενός δικτύου MLP με δύο κρυφά επίπεδα (οι συνδέσεις των νευρώνων που παραλείπονται θεωρούνται δεδομένες ενώ επίσης εννοείται και η ύπαρξη νευρώνων πόλωσης)



Τα κρυφά επίπεδα του δικτύου, κατά την διαδικασία εκπαίδευσης, προσαρμόζονται ώστε να εξάγουν τα κυριότερα χαρακτηριστικά των δεδομένων εισόδου. Η λειτουργία αυτή επιτελείται μέσω ενός μη γραμμικού μετασχηματισμού από τα δεδομένα εισόδου σε έναν νέο χώρο, τον χώρο χαρακτηριστικών (feature space) στον οποίο τα δεδομένα του εκάστοτε προβλήματος είναι πιο εύκολα διαχωρίσιμα [42]. Σε αντίθεση όμως με το επίπεδο εξόδου, για τα κρυμμένα επίπεδα δεν υπάρχει κάποια συγκεκριμένη επιθυμητή απόκριση. Συνεπώς, δημιουργείται το πρόβλημα του πως πρέπει να προσαρμόζονται τα βάρη τους για να ελαχιστοποιηθεί το σφάλμα στην έξοδο. Αυτό το πρόβλημα μπορεί να αντιμετωπιστεί με την χρήση μίας, μη γραμμικής, διαφορίσιμης συνάρτησης ενεργοποίησης, αντί για την συνάρτηση προσήμου που χρησιμοποιούταν στο μοντέλο Perceptron, η οποία επιτρέπει την συνεχή βελτιστοποίηση. Έτσι οδηγούμαστε στην γενίκευση του delta rule για δίκτυα πολλαπλών στρωμάτων, τον οπισθοβατικό αλγόριθμο (backpropagation algorithm) ο οποίος περιγράφεται στην επόμενη υποενότητα.

2.1.3 Βασικές αρχές εκπαίδευσης Νευρωνικών Δικτύων

Στην απλούστερη περίπτωση του νευρωνικού δικτύου με ένα επίπεδο, η διαδικασία εκπαίδευσης είναι σχετικά απλή λόγω του γεγονότος πως το σφάλμα (συνάρτηση απώλειας ή loss function) μπορεί να γραφεί ως άμεση συνάρτηση των βαρών. Για ένα δίκτυο όμως με πολλαπλά επίπεδα το αντίστοιχο σφάλμα είναι μια αρκετά πιο σύνθετη συνάρτηση σύνθεσης του συνόλου των βαρών του μοντέλου. Σε αυτήν την περίπτωση για τον υπολογισμό του δυναμικού (gradient) της συνάρτησης σφάλματος που απαιτείται για την ανανέωση των βαρών του δικτύου χρησιμοποιείται ο οπισθοβατικός αλγόριθμος. Ο οπισθοβατικός αλγόριθμος έχει συνεισφέρει σε μεγάλο βαθμό στην πολύ ευρεία χρήση των νευρωνικών δικτύων, καθώς αυτός αφαιρεί από τον ερευνητή το βάρος της άμεσης επεξεργασίας όλων των βημάτων ανανέωσης των βαρών, κάτι που καθίσταται απαγορευτικά σύνθετο αχόμη και για σχετικά απλές αρχιτεκτονικές [1].

Πιο αναλυτικά, ο αλγόριθμος αυτός περιλαμβάνει δύο φάσεις, τη φάση προς τα εμπρός και τη φάση προς τα πίσω. Στην πρώτη, προκαταρτική, φάση τροφοδοτούνται στο δίκτυο οι είσοδοι του συνόλου δεδομένων εκπαίδευσης και υπολογίζονται οι έξοδοι όλων των νευρώνων. Στην δεύτερη φάση υπολογίζει το δυναμικό της συνάρτησης σφάλματος ως προς τα βάρη του νευρωνικού δικτύου, χρησιμοποιώντας τον διαφορικό κανόνα της αλυσίδας ($\frac{dz}{dx} = \frac{dz}{dy} \cdot \frac{dy}{dx}$). Οι υπολογισμοί αυτοί γίνονται για ένα επίπεδο τη φορά, ξεκινώντας από το επίπεδο εξόδου και προχωρώντας προς χαμηλότερα επίπεδα, ώστε να αποφεύγεται η εκτέλεση των υπολογισμών για τις ενδιάμεσες παραγώγους του κανόνα της αλυσίδας πολλαπλές φορές. Χρησιμοποιώντας μαθηματικούς συμβολισμούς, η ανανέωση των βαρών με τη μέθοδο της κατάβασης δυναμικού γίνεται ως εξής:

$$\Delta w_{ij} = -a \cdot \frac{\partial E}{\partial w_{ij}}$$

όπου το α και το E, όπως και στην απλή περίπτωση του Perceptron, συμβολίζουν τον ρυθμό μάθησης και τη συνάρτηση απώλειας αντίστοιχα ενώ το w_{ij} είναι το βάρος της ακμής από τον κόμβο j στον κόμβο i.

Εφαρμόζοντας τον κανόνα της αλυσίδας θα ισχύει:

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_i} \cdot \frac{\partial y_i}{\partial v_i} \cdot \frac{\partial v_i}{\partial w_{ij}}$$

όπου v_i και y_i είναι οι έξοδοι του νευρώνα
ἰπριν και μετά την εφαρμογή της συνάρτησης ενεργοποίησης f και συνεπώς προκύπτει:

$$\frac{\partial E}{\partial w_{ij}} = \frac{\partial E}{\partial y_i} \cdot f'(v_i) \cdot v_j$$

Αχολούθως, αν ο νευρώνας i είναι νευρώνας εξόδου η ποσότητα $\frac{\partial E}{\partial y_i}$ θα προχύπτει ως η μεριχή παράγωγος της συνάρτησης σφάλματος ως προς το σφάλμα εξόδου. Στην περίπτωση που ως συνάρτηση απώλειας χρησιμοποιείται το μέσο τετραγωνιχό σφάλμα, αυτή θα ισούται με $(d_i - y_i)$, όπου d_i είναι η επιθυμητή έξοδος.

Αν ο νευρώνας i δεν ανήκει στο επίπεδο εξόδου αλλά σε κάποιο άλλο επίπεδο l, τότε από τον κανόνα αθροίσματος της παραγώγου θα ισχύει:

$$\frac{\partial E}{\partial y_i} = \sum_{k \in (l+1)} w_{ki} \cdot \frac{\partial E}{\partial y_k} \cdot f'(v_k)$$

όπου οι ποσότητες $\partial E/\partial y_k$, $f'(v_k)$ θα είναι ήδη γνωστές από την εκτέλεση του οπισθοβατικό αλγορίθμου για το επόμενο επίπεδο. Εδώ αξίζει να σημειωθεί πως για την ανανέωση των βαρών ενός νευρωνικού δικτύου μπορούν να χρησιμοποιηθούν και άλλες μέθοδοι ελαχιστοποίησης πέραν της απλής κατάβασης δυναμικού όμως η γενική λογική του οπισθοβατικό αλγορίθμου εξακολουθεί να ισχύει.

2.1.4 Συναρτήσεις ενεργοποίησης Νευρωνικών Δικτύων

Η συνάρτηση ενεργοποίησης ενός νευρώνα εφαρμόζεται στο άθροισμα των εισόδων του για να παράξει την τελική έξοδο του. Συνεπώς η επιλογή της έχει σημαντική επίπτωση στην συμπεριφορά ενός νευρωνικού δικτύου. Υπενθυμίζεται ότι στην περίπτωση του Perceptron είχε χρησιμοποιηθεί για αυτό το σκοπό η συνάρτηση προσήμου με αποτέλεσμα αυτό να μπορεί να χρησιμοποιηθεί μόνο για προβλήματα δυαδικής ταξινόμησης. Παράλληλα η απλή αρχιτεκτονική του το περιόριζε στην επίλυση μόνο γραμμικώς διαχωρίσιμων προβλημάτων, για την αντιμετώπιση των οποίων προτάθηκε το μοντέλο MLP με την προσθήκη κρυφών επιπέδων τα οποία μετασχηματίζουν τα δεδομένα εισόδου σε έναν χώρο όπου αυτά μπορούν να διαχωριστούν γραμμικά. Η προαναφερθείσα όμως η λειτουργία είναι δυνατή μόνο με την απαίτηση του αλγορίθμου οπισθοδρόμησης για διαφορισιμότητα, οδήγησαν στην υιοθέτηση συναρτήσεων όπως η σιγμοειδής και η tanh, οι οποίες για πολλά χρόνια ήταν οι κατεξοχήν ευρέως χρησιμοποιούμενες μη γραμμικές συναρτήσεις ενεργοποίησης νευρωνικών δικτύων.

Η περιορισμένη όμως, λόγω κορεσμού, ταχύτητα αυτών των συναρτήσεων μπορεί να βελτιωθεί σημαντικά με την χρήση της μη γραμμικής συνάρτησης ReLU (f(x) = max(x,0)) η οποία δεν παρουσιάζει φαινόμενα κορεσμού [11]. Η διαφορά αυτή

είναι αχόμη πιο αισθητή για ιδιαίτερα σύνθετα νευρωνιχά δίχτυα. Ως αποτέλεσμα, οι χλασιχές συναρτήσεις ενεργοποίησης που υπόχεινται σε χορεσμό έχουν πλέον αρχετά περιορισμένη χρήση, πέραν του επιπέδου εξόδου ενός διχτύου. Εδώ να σημειώθεί πως επειδή η ReLU είναι διαφορίσιμη για όλες τις πραγματιχές εισόδους εχτός από ένα μόνο σημείο, αυτή είναι συμβατή με τον οπισθοβατιχό αλγόριθμο.



Σχήμα 2.1.4: Κοινές συναρτήσεις ενεργοποίησης νευρωνικών δικτύων

2.1.5 Συνελικτικά Νευρωνικά Δίκτυα

Τα συνελικτικά νευρωνικά δίκτυα (convolutional neural networks ή CNN) είναι ένας ειδικός τύπος βιολογικά εμπνευσμένων νευρωνικών δικτύων χωρίς ανάδραση που χρησιμοποιούν συνελικτικά φίλτρα αντί για πλήρως συνδεδεμένους νευρώνες σε τουλάχιστον ένα επίπεδο τους. Αυτά τα μοντέλα έχουν εφαρμογή σε δεδομένα με γνωστή πλεγματική δομή που παρουσιάζουν χωρική συσχέτιση, όπως εικόνες και βίντεο, καθώς ενσωματώνουν την υπόθεση της χωρικής συσχέτισης στην ίδια τους την αρχιτεκτονική. Τα δίκτυα αυτά παίρνουν το όνομα τους από την μαθηματική πράξη της συνέλιξης, η οποία σε αυτήν την περίπτωση χρησιμοποιείται για την χαρτογράφηση διαφόρων χαρακτηριστικών ενός επιπέδου σε ένα άλλο. Συγκεκριμένα η πράξη της συνέλιξης στη μία διάσταση για δύο αχολουθίες f, g ορίζεται ως εξής:

$$(f * g)[n] = \sum_{m = -\infty}^{\infty} f[n]g[n - m]$$

Η συνέλιξη όμως ορίζεται και χρησιμοποιείται και για περισσότερες διαστάσεις. Ειδικότερα, η συνέλιξη στις δύο διαστάσεις, που χρησιμοποιείται ευρέως σε εφαρμογές όρασης υπολογιστών, ορίζεται ως εξής:

$$(f*g)[i][j] = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} f[n][m] \cdot g[i-n][j-m] = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} f[i-n][j-m] \cdot g[n][m]$$

Έτσι πραχτιχά η συνέλιξη δύο αχολουθιών υπολογίζεται αναστρέφοντας τη μια χαι υπολογίζοντας το εσωτεριχό γινόμενο αυτής με την άλλη για διαφορετιχές σχετιχές μετατοπίσεις των δύο, με βήμα 1. Πολλές όμως υλοποιήσεις συνελιχτιχών νευρωνιχών διχτύων χρησιμοποιούν μια συνάρτηση που προσομοιάζει την συνέλιξη αλλά χωρίς την αναστροφή του φίλτρου, η οποία ονομάζεται ετεροσυσχέτιση (cross-correlation) χαι ορίζεται στις δύο διαστάσεις από τον παραχάτω τύπο [6]:

$$(f*g)[i][j] = \sum_{n=-\infty}^{\infty} \sum_{m=-\infty}^{\infty} f[i+n][j+m] \cdot g[n][m]$$

Χάριν ευχολίας, σε αυτήν την διπλωματιχή εργασία θα αναφέρομαι στο εξής τόσο στην συνέλιξη χαι την ετεροσυσχέτιση όσο χαι στις παραλλαγές τους με βήμα μεγαλύτερο της μονάδας ως συνέλιξη. Εδώ να σημειώσω επίσης πως η συνέλιξη δύο πεπερασμένων αχολουθιών μπορεί να υπολογιστεί με την διαδιχασία που περιγράφηχε παραπάνω απλώς θεωρώντας πως αυτές έχουν μηδενιχές τιμές για δείχτες μεγαλύτερους των διαστάσεων τους. Στις εφαρμογές των CNN η είσοδος χάθε συνελιχτιχού επιπέδου είναι συνήθως ένας πολυδιάστατος πίναχας (tensor) δεδομένων ενώ το φίλτρο είναι ένας πολυδιάστατος (τρισδιάστατος για τις περισσότερες εφαρμογές) πίναχας με το ίδιο βάθος με το τρέχον επίπεδο εισόδου με διαστάσεις $H_1 \cdot W_1 \cdot D_1$ (το οποίο ισοδυναμεί με μία ειχόνα $H_1 \cdot W_1$ με D_1 χανάλια), Κ φίλτρα $F \cdot F \cdot D_1$ (το οποίο ισοδυναμεί με Κ σύνολα δισδιάστατων συνελιχτιχών φίλτρων D_1 χαναλιών), βήμα (stride) S και padding (αύξηση του ύψους και πλάτους της εισόδου ψε πιπέδου θα έχει διαστάσεις:

$$H_2 \cdot W_2 \cdot D_2$$
$$W_2 = \frac{(W_1 - F + 2P)}{S}$$
$$H_2 = \frac{(H_1 - F + 2P)}{S}$$



Σχήμα 2.1.5: Παράδειγμα δισδιάστατης συνέλιξης δυο πινάχων

 $D_2 = K$

Σε ένα συνελιχτικό δίκτυο οι συνδέσεις είναι ιδιαίτερα αραιές (sparse connectivity), καθώς οι τιμές εξόδου ενός συνελικτικού επιπέδου προκύπτουν από το εσωτεριχό γινόμενο μίας μιχρής χωριχής περιοχής της εισόδου και ενός κοινού συνελιπτιπού φίλτρου (parameter sharing). Έτσι ένα τέτοιο μοντέλο μπορεί να εντοπίζει συγκεκριμένα χαρακτηριστικά σε διαφορετικά σημεία της εισόδου. Το γεγονός αυτό επίσης συνεπάγεται πως τα CNN παρουσιάζουν σημαντική μεταθετική αμεταβλητότητα (translation invariance) και επιτρέπουν την εξαγωγή από κάθε επίπεδο χωριχών χαραχτηριστιχών της εισόδου, από τα πιο απλά, όπως οριζόντιες χαι χάθετες γραμμές, στα πιο σύνθετα χαθώς ανεβαίνουμε σε πιο υψηλά επίπεδα. Αχόμη οι αραιές συνδέσεις των συνελιχτικών διχτύων συνεπάγονται μικρό αριθμό παραμέτρων σε σχέση με ένα πλήρως συνδεδεμένο δίκτυο για μία παρόμοια εφαρμογή το οποίο οδηγεί σε ταχύτερη εκπαίδευση και μικρότερες απαιτήσεις σε μνήμη. Παράλληλα, ο μικρός αριθμός παραμέτρων κάνει τα CNN λιγότερο επιρρεπή στο πρόβλημα του overfitting, όπου ένα μοντέλο μηχανικής μάθησης παρουσιάζει σημαντικά υψηλότερη αχρίβεια στο σύνολο δεδομένων εχπαίδευσης απ' ό,τι σε αυτό της αξιολόγησης. Αυτό συμβαίνει λόγω του ότι το μοντέλο έχει μάθει να αναγνωρίζει τα δεδομένα εκπαίδευσης βασιζόμενο σε ασήμαντες (για το εν λόγω πρόβλημα) λεπτομέρειες τους.

Σε ένα τυπικό συνελικτικό επίπεδο μετά τον υπολογισμό της συνέλιξης εφαρμόζεται στην έξοδο του κάποια (μη γραμμική) συνάρτηση ενεργοποίησης όπως η ReLU. Συχνά αυτό το βήμα ακολουθείται από κάποιο επίπεδο ομαδοποίησης (pooling layer) στο οποίο η έξοδος υποδειγματοληπτείται, γεγονός που βελτιώνει περαιτέρω την μεταθετική αμεταβλητότητα του μοντέλου. Δύο από τα πιο συχνά χρησιμοποιούμενα τέτοια επίπεδα είναι αυτά του max pooling και average pooling. Στην πρώτη περίπτωση για χάθε μικρή περιοχή (π.χ. 2x2) χάθε επιπέδου της εξόδου επιστρέφεται η μέγιστη τιμή εντός αυτής, ενώ στην δεύτερη για χάθε τέτοια περιοχή επιστρέφεται ο μέσος όρος. Αχόμη, σε πολλές αρχιτεκτονιχές CNN τα τελευταία επίπεδα υλοποιούνται ως πλήρως συνδεδεμένα στρώματα για να αυξηθεί η ισχύς των υπολογισμών προς το τέλος [1].

Σχήμα 2.1.6: LeNet 5: Ένα από τα πρώτα CNN που αποτελεί χαρακτηριστικό παράδειγμα της τυπικής αρχιτεκτονικής τους



2.2 MOS τρανζίστορ

To MOS (metal oxide semiconductor) τρανζίστορ, η αλλιώς MOSFET (metal oxide semiconductor field-effect transistor), αποτελεί την βασική κυκλωματική δομική μονάδα των περισσοτέρων σύγχρονων ψηφιακών αλλά και αναλογικών ολοκληρωμένων κυκλωμάτων. Σήμερα, σε σύγκριση με τα BJT (bipolar junction transistor), τον άλλον ευρέως χρησιμοποιούμενο τύπο τρανζίστορ, αυτά μπορούν να κατασκευαστούν πιο εύκολα και σε πολύ μικρότερες διαστάσεις ενώ τυπικά απαιτούν αρκετά λιγότερη ισχύ [45]. Παράλληλα τα MOSFET χαρακτηρίζονται από μεγάλη ευελιξία στην λειτουργία τους, καθώς μπορούν να χρησιμοποιηθούν στη θέση πολλών άλλων κλασικών κυκλωμάτων όπως πυκνωτές και αντιστάσεις. Έτσι, επιτρέπουν την παραγωγή αναλογικών και ψηφιακών κυκλωμάτων πολύ υψηλής ολοκλήρωσης που χρησιμοποιούν σχεδόν αποκλειστικά τέτοια ημιαγώγιμα στοιχεία [47].

2.2.1 Βασική δομή MOS τρανζίστορ

Στο Σχήμα 2.2.1 παρουσιάζεται η βασική δομή ενός MOS τρανζίστορ τύπου n (NMOS). Ένα τέτοιο στοιχείο κατασκευάζεται επάνω σε ένα δισκίο μονοκρυσταλ-

λικού πυριτίου στο οποίο έχει δημιουργηθεί ένα υπόστρωμα τύπου p (ημιαγωγού με προσμίξεις αποδεκτών ηλεκτρονίων). Σε αυτό το υπόστρωμα δημιουργούνται δύο περιοχές τύπου n⁺ (ημιαγωγού με ισχυρή νόθευση δοτών ηλεκτρονίων) οι οποίες λειτουργούν ως περιοχές πηγής (source) και υποδοχής (drain). Ακολούθως, εναποτίθεται πάνω από το υπόστρωμα, ανάμεσα από τις παραπάνω δύο περιοχές, ένα λεπτό επίπεδο από (συνήθως) διοξειδιο του πυριτίου (SiO₂), που λειτουργεί ως μονωτής. Επάνω σε αυτό εναποτίθεται μέταλλο ή ισχυρά νοθευμένο πολυκρυσταλλικό πυρίτιο (που λειτουργεί ως καλός αγωγός) [46], το οποίο αποτελεί το ηλεκτρόδιο πύλης (gate). Επίσης τοποθετούνται αγώγιμες επαφές στην πηγή, την υποδοχή και το υπόστρωμα (που επίσης αποκαλείται σώμα η body) και έτσι δημιουργούνται οι τέσσερεις ακροδέκτες του MOSFET: G, D, S και B.

Για το PMOS τρανζίστορ η δομή είναι παρόμοια αλλά αντί για n^+ ημιαγωγοί για την πηγή και την υποδοχή χρησιμοποιούνται ημιαγωγοί τύπου p^+ ενώ το υπόστρωμα κατασκευάζεται από ημιαγωγό τύπου n.

2.2.2 Βασική λειτουργία μεγάλου σήματος MOS τρανζίστορ

Η παραχάτω ανάλυση γίνεται για ένα MOSFET n χαναλιού, όμως η ανάλυση για ένα PMOS τρανζίστορ είναι ανάλογη.

Αν σε ένα NMOS τρανζίστορ εφαρμοστεί μια θετική τάση (αρνητική για PMOS) V_{GS} μεταξύ της πηγής και της πύλης, οι οπές του υποστρώματος κάτω από την πύλη απωθούνται δημιουργώντας μια περιοχή απογύμνωσης. Παράλληλα, η πύλη έλχει ηλεχτρόνια τα οποία με την επαρχή συσσώρευση τους δημιουργούν ένα χανάλι αγωγής ρεύματος τύπου n μεταξύ της πηγής και της υποδοχής (εξ' ου και το όνομα MOSFET n καναλιού). Έτσι ουσιαστικά σχηματίζεται ένας πυκνωτής μεταξύ της πύλης και του καναλιού αγωγής, όπου το μονωτικό στρώμα δρα ως διηλεκτρικό, χαι το εγχάρσιο ηλεχτρικό πεδίο που δημιουργείται ελέγχει την αγωγιμότητα του ${\rm n}$ χαναλιού [45]. Η τιμή της τάσης V_{GS} για την οποία συσσωρεύεται επαρχής αριθμός ελεύθερων ηλεκτρονίων μεταξύ υποδοχής και πύλης για να σχηματιστεί κανάλι αγωγής ονομάζεται τάση κατωφλίου (threshold voltage) και συχνά συμβολίζεται V_t . Το φορτίο στο κανάλι καθορίζεται από την τάση $V_{ov} = V_{GS} - V_t$ που ονομάζεται τάση υπεροδήγησης (overdrive voltage). Αν η τάση V_{GS} είναι μικρότερη από την τάση Vt θεωρώ αρχικά για απλότητα ότι το τρανζίστορ δεν άγει καθόλου και έτσι βρίσκεται στην περιοχή αποκοπής. Στην πραγματικότητα, για τάσεις V_{GS} λίγο χαμηλότερες από την τάση κατωφλίου, το NMOS βρίσκεται στην περιοχή υποκατωφλίου (subthreshold region) και άγει πολύ μικρά ρεύματα. Για θετικές τιμές της τάσης υπεροδήγησης, ανάλογα με την τιμή της τάσης μεταξύ της υποδοχής και της πηγής V_{DS} το NMOS τρανζίστορ θα βρίσκεται σε μία εκ των δύο παρακάτω περιοχών:

1. Για $V_{DS} < V_{ov}$ το MOSFET
n καναλιού θα βρίσκεται στην περιοχή τριόδου. Σε αυτήν την περιοχή το ρεύμα που άγει το κανάλι εξαρτάται άμεσα από την

τιμή της V_{DS} σύμφωνα με την παρακάτω σχέση:

$$I_D=k_n'(\frac{W}{L})[V_{ov}V_{DS}-\frac{1}{2}V_{DS}^2]$$

Ο συντελεστής $k'_n = m_n C_{ox}$ ονομάζεται παράμετρος διαγωγιμότητας τεχνολογίας κατασκευής (process transconductance) και είναι ίσος με το γινόμενο της κινητικότητας των ηλεκτρονίων (m_n) και της χωρητικότητας οξειδίου (C_{ox}) , ενώ τα W και L είναι το πλάτος και το μήκος του καναλιού αντίστοιχα. Επίσης, από τον παραπάνω τύπο φαίνεται πως για μικρές V_{DS} η εξάρτηση του ρεύματος υποδοχής από αυτό είναι γραμμική, εφόσον ο τετραγωνικός όρος θα είναι αμελητέος, και άρα το τρανζίστορ λειτουργεί ως ωμική αντίσταση.

2. Στην περιοχή λειτουργίας που περιγράφηκε παραπάνω, όσο αυξάνεται η V_{DS} το βάθος του καναλιού στην πλευρά της υποδοχής γίνεται όλο και μικρότερο. Όταν όμως η V_{DS} φτάσει την τάση κατωφλίου, υφίσταται στραγγαλισμός του καναλιού (pinched off channel). Σε αυτό το σημείο το NMOS εισέρχεται στην περιοχή κορεσμού και η περαιτέρω αύξηση αυτής της τάσης ($V_{DS} \ge V_{ov}$) έχει πολύ μικρή επίδραση στο σχήμα και στο φορτίο του καναλιού [45]. Σε αυτήν την περιοχή το ρεύμα που άγει το κανάλι δίνεται προσεγγιστικά από τον παρακάτω τύπο:

$$I_D = \frac{1}{2}k'_n(\frac{W}{L})V_{ov}^2$$

Παρά όμως την προηγούμενη απλοϊκή προσέγγιση, η αύξηση της V_{DS} έχει επίδραση στο κανάλι αγωγής. Συγκεκριμένα, καθώς αυτή αυξάνεται, το σημείο στραγγαλισμού μετατοπίζεται ελαφρώς από την πηγή προς την υποδοχή (διαμόρφωση μήκους καναλιού ή φαινόμενο Early). Αυτό το φαινόμενο μπορεί να ενσωματωθεί στην προηγούμενη σχέση του ρεύματος καναλιού (σε κατάσταση κορεσμού) ως εξής:

$$I_D = \frac{1}{2}k'_n(\frac{W}{L})V_{ov}^2(1+\lambda V_{DS})$$

όπου το λ είναι μια παράμετρος που εξαρτάται από την τεχνολογία κατασκευής και το μήκος καναλιού.

Εδώ αξίζει να σημειώσω πως η τάση κατωφλίου αν και μπορεί να θεωρηθεί προσεγγιστικά σταθερή, παρουσιάζει μη αμελητέα εξάρτηση από την τάση V_{SB} , σύμφωνα με τον τύπο:

$$V_t = V_{t0} + \gamma [\sqrt{2\phi_f + V_{SB}} - \sqrt{2\phi_f}]$$

Όπου V_{t0} είναι η τάση κατωφλίου για $V_{SB} = 0$, το ϕ_f είναι μια φυσική παράμετρος και το γ είναι μια παράμετρος εξαρτώμενη από την τεχνολογία κατασκευής.

2.2.3 Λειτουργία MOS τρανζίστορ στην περιοχή υποκατωφλίου

Στην παραπάνω ανάλυση έγινε η υπόθεση πως για τάσεις V_{GS} μικρότερες από την τάση κατωφλίου το MOSFET μπαίνει απότομα σε αποκοπή. Όμως παρά αυτήν την απλοϊκή προσέγγιση, στην πραγματικότητα για τέτοιες τάσεις V_{GS} κοντά στην V_t συνεχίζει να υπάρχει ένα «ασθενές» στρώμα αναστροφής (weak inversion). Έτσι για V_{GS} λίγο μικρότερη από την τάση κατωφλίου και τάση V_{DS} μεγαλύτερη από περίπου 100mV υπάρχει ένα πεπερασμένο ρεύμα που εξαρτάται εκθετικά από την V_{GS} ως εξής:

$$I_D = I_0 \cdot exp(\frac{V_{GS}}{\xi V_\tau})$$

Όπου το I_0 είναι ανάλογο του λόγου W/L, το $\xi > 1$ είναι ένας συντελεστής μη ιδανικότητας και V_{τ} είναι η θερμική τάση που ισούται με $\frac{kT}{q}$ (το k συμβολίζει την σταθερά Boltzman, το T την απόλυτη θερμοκρασία και το q είναι το απόλυτο φορτίο του ηλεκτρονίου) [46].

Αυτή η μικρή αγωγή ρεύματος κάτω από την τάση κατωφλίου συχνά αποτελεί σημαντικό πρόβλημα για κυκλώματα με πολύ μεγάλο αριθμό τρανζίστορ. Αυτό συμβαίνει, καθώς ενώ για ένα MOSFET το σχετικό ρεύμα μπορεί να είναι αμελητέο, το αθροιστικό ρεύμα εκατομμυρίων τέτοιων στοιχείων μπορεί να προκαλεί σημαντική στατική κατανάλωση ισχύος, ενώ σε κυκλώματα μνήμης μπορεί να οδηγήσει σε απώλεια αναλογικής πληροφορίας [45]. Παράλληλα όμως, η λειτουργία των MOS τρανζίστορ στην περιοχή υποκατωφλίου επιτρέπει την σχεδίαση αναλογικών κυκλωμάτων εξαιρετικά χαμηλής κατανάλωσης ισχύος με κόστος την ταχύτητα λειτουργίας.



Σχήμα 2.2.1: Φυσική δομή NMOS τρανζίστορ: Στην επάνω εικόνα παρουσιάζεται μια προοπτική άποψη και στην κάτω μια τομή

Σχήμα 2.2.2: Τομή ολοκληρωμένου CMOS (complementary MOS) κυκλώματος όπου για το PMOS τρανζίστορ κατασκευάζεται μια ξεχωριστή περιοχή n ημιαγωγού (n well) στο p υπόστρωμα



Σχήμα 2.2.3: Μορφή καναλιού NMOS για τρείς περιπτώσεις: (επάνω) $V_{DS} < V_{ov},$ (μέση) $V_{DS} = V_{ov},$ (κάτω) $V_{DS} > V_{ov}$





Σχήμα 2.2.4: Βασικές περιοχές λειτουργίας NMOS τρανζίστορ

Σχήμα 2.2.5: Εξάρτηση I_D NMOS τρανζίστορ από τη
ν V_{GS} στις περιοχές υποκατωφλίου και κορεσμού



Κεφάλαιο 3

Δομικά Κυκλώματα

Σε αυτό το κεφάλαιο αναλύεται ο σχεδιασμός των βασικών κυκλωματικών δομικών μονάδων που χρησιμοποιούνται για την υλοποίηση του εν λόγω αναλογικού συνελικτικού νευρωνικού δικτύου. Συγκεκριμένα, παρουσιάζεται ο σχεδιασμός και οι βασικές αρχές λειτουργίας των κυκλωμάτων των πολλαπλασιαστών, της συνάρτησης ενεργοποίησης ReLU και του τελεστή argmax.

3.1 Κυχλώματα υλοποίησης σιγμοειδούς συνάρτησης

Στον πυρήνα των χυχλωμάτων αναλογικών πολλαπλασιαστών υψηλής αχρίβειας και εξαιρετικά χαμηλής κατανάλωσης που χρησιμοποιούνται στην προτεινόμενη υλοποίηση του αναλογικού νευρωνικού δικτύου, βρίσκεται ένα κύκλωμα σιγμοειδούς. Συγκεκριμένα, το κύκλωμα σιγμοειδούς αναφέρεται σε ένα κύκλωμα το οποίο δύναται να παράγει έξοδο ρεύματος της οποίας η εξάρτηση από μία τάση εισόδου προσομοιάζει την σιγμοειδή συνάρτηση $(1/(1 + e^{-x}))$.

3.1.1 Απλό χύχλωμα σιγμοειδούς με διαφοριχό ζεύγος

Μια από τις απλούστερες δυνατές υλοποιήσεις ενός τέτοιου κυκλώματος είναι με την χρήση ενός απλού διαφορικού ζεύγους, όπως φαίνεται στο Σχήμα 3.1.1¹. Συγκεκριμένα, το εν λόγω κύκλωμα περιλαμβάνει αρχικά ένα διαφορικό ζεύγος το οποίο τροφοδοτείται με ρεύμα πόλωσης μέσω ενός απλού NMOS καθρέπτη ρεύματος. Ακολούθως, κάθε ένα από τα M1, M2 τροφοδοτεί ένα διοδικά συνδεδεμένο PMOS τρανζίστορ, το δεξιά εκ των οποίων σχηματίζει έναν απλό PMOS καθρέπτη ρεύματος για να απομονωθεί η τελική έξοδος, όπως φαίνεται και στο σχετικό σχεδιάγραμμα.

¹Στα χυχλωματικά διαγράμματα της παρούσας διπλωματικής, σε όσα τρανζίστορ δεν εμφανίζεται ο αχροδέχτης σώματος εννοείται πως αυτός είναι συνδεδεμένος με την αντίστοιχη τάση τροφοδοσίας (0.3V για PMOS και -0.3V για NMOS)
Προσομοιώνοντας αυτό το χύχλωμα για Vr = 0, Ibias = 10nA και μεταβλητή Vin, προχύπτει η έξοδος του Σχήματος 3.1.2 η οποία προσεγγίζει μια ανεστραμμένη χαι μετατοπισμένη σιγμοειδή συνάρτηση. Η μορφή αυτή της εξόδου προχύπτει, χαθώς στην αρχική κατάσταση, όπου η Vin είναι ίση με την τάση αρνητικής τροφοδοσίας (-300mV), η τάση V_{GS} του M1 είναι πολύ μικρότερη από αυτήν του M2 συνεπώς το δεύτερο άγει σχεδόν όλο το ρεύμα πόλωσης που παρέχει ο χαθρέπτης ρεύματος. Καθώς αυξάνεται η V_{GS} του M1, το ρεύμα που αυτό άγει αυξάνεται αρχικά με προσεγγιστικά εκθετικό ρυθμό, ο οποίος όμως στην συνέχεια μειώνεται όσο το ρεύμα αυτό πλησιάζει το ρεύμα πόλωσης του διαφορικού ζεύγους, έως ότου τελικά φτάσει αυτήν την τιμή για τάσεις V_{GS} κοντά στην θετική τροφοδοσία. Παράλληλα με το Μ1, το ρεύμα που άγει το Μ2, από το οποίο προχύπτει χαι η έξοδος σε αυτήν την περίπτωση, θα αχολουθεί συμπληρωματιχή πορεία δεδομένου ότι το συνολιχό ρεύμα πόλωσης που παρέχει ο καθρέπτης στο διαφορικό ζεύγος είναι σχεδόν σταθερό. Οι σιγμοειδείς καμπύλες αυτές μπορούν επίσης να μετατοπίζονται με την χρήση διαφορετικών τιμών για την Vr. Αυτή όμως είναι μια ιδιότητα που δεν χρησιμοποιείται στην παρούσα σχεδίαση και έτσι στην παρακάτω ανάλυση θεωρώ ότι η τάση Vr είναι σταθερή και ίση με 0V.







Σχήμα 3.1.2: Παράδειγμα εξόδου απλού κυκλώματος σιγμοειδούς για $\mathrm{Vr}=0,$ Ibias $=10\mathrm{nA}$

Το χύχλωμα όμως αυτό παρουσιάζει αχόμη μια ιδιότητα που είναι ιδιαίτερα χρήσιμη για την παρούσα εφαρμογή, η οποία φαίνεται χαθαρά στο Σχήμα 3.1.3. Συγκεχριμένα, η έξοδος ρεύματος αυτού έχει σε μεγάλο βαθμό γραμμιχή εξάρτηση από την τιμή του Ibias, για σχεδόν όλο το εύρος της τάσης Vin. Έτσι η συνάρτηση μεταφοράς (H(Vin) = Iout/Ibias) του παραπάνω χυχλώματος είναι σχετιχά σταθερή για διαφορετιχές τιμές του ρεύματος Ibias (Σχήμα 3.1.4). Από την συνάρτηση αυτή μεταφοράς φαίνεται πως η έξοδος του εν λόγω χυχλώματος είναι προσεγγιστιχά ίση με το ρεύμα Ibias επί έναν συντελεστή που χυμαίνεται μεταξύ μιας ελάχιστης τιμής χοντά στο 0 χαι μίας μέγιστης τιμής (≈ 2 για το συγχεχριμένο παράδειγμα) που εξαρτάται από τις διαστάσεις των χρησιμοποιούμενων τρανζίστορ. Συνεπώς μπορώ να χρησιμοποιήσω ένα τέτοιο χύχλωμα ως έναν πολύ απλό αναλογικό πολλαπλασιαστή ενός ρεύματος με έναν θετιχό συντελεστή που χυμαίνεται εντός χάποιου εύρους τιμών.

Όπως φαίνεται όμως και από τα Σχήματα 3.1.3 και 3.1.4, το παραπάνω κύκλωμα ως αναλογικός πολλαπλασιαστής αντιμετωπίζει κάποια σημαντικά προβλήματα. Πρώτον και βασικότερο είναι το γεγονός πως ενώ η σχέση μεταξύ του Ιουτ και του Ibias για δεδομένη Vin φαίνεται σε μεγάλο βαθμό γραμμική, η ευθεία που ορίζει αυτήν την σχέση δεν διέρχεται ιδιαίτερα κοντά από το σημείο (0, 0) όπως και είναι επιθυμητό για μια τέτοια εφαρμογή. Δεύτερον, η παραπάνω σχέση για χαμηλές τιμές Ibias και ιδιαίτερα για πολύ αρνητικές τιμές Vin (που συνεπάγονται μεγάλους



Σχήμα 3.1.3: Παράδειγμα εξόδου απλού κυκλώματος σιγμο
ειδούς συναρτήσει του Ibias για $\mathrm{Vr}=0$

συντελεστές πολλαπλασιασμού) παύει να είναι γραμμική, όπως φαίνεται στο Σχήμα 3.1.3 όπου παρατηρείται μια (ήπια) καμπυλότητα.

3.1.2 Κύκλωμα σιγμοειδούς με διαφορικό ζεύγος διαφοράς

Μια πρώτη αλλά σημαντική βελτίωση στο παραπάνω απλό κύκλωμα ενός αναλογικού πολλαπλασιαστή μπορεί να γίνει με την χρήση ενός διαφοριχού ζεύγους διαφοράς (differential difference pair) αντί για ένα απλό διαφορικό ζεύγος (Σχήμα 3.1.5). Το διαφορικό ζεύγος διαφοράς αποτελείται από δύο επιμέρους διαφορικά ζεύγη που χάθε ένα τους παράγει σιγμοειδή έξοδο με ρυθμιζόμενη χλίση [48]. Αχόμη, στους αχροδέχτες σώματος των M01, M04 εφαρμόζεται τάση Vc η οποία μέσω του φαινομένου σώματος μπορεί να μεταβάλλει αισθητά τη μορφή της εξόδου ρεύματος (Σχήμα 3.1.6). Για την τροφοδοσία του ρεύματος πόλωσης στο διαφορικό ζεύγος διαφοράς χρησιμοποιείται κασκοδικός καθρέπτης ρεύματος, αντί για απλό όπως στο αρχικό χύχλωμα σιγμοειδούς, ώστε αυτό να παρουσιάζει σχεδόν αμελητέα μεταβολή για διαφορετικές τιμές τάσεων εισόδου. Επίσης, για την απομόνωση της εξόδου, παρατηρώ μέσω προσομοίωσης του χυχλώματος πως ο διπλός χαθρέπτης ρεύματος δεν προσφέρει χάποιο ουσιαστιχό πλεονέχτημα χαι έτσι επιλέγεται για αυτό το σχοπό ένας απλός καθρέπτης. Τέλος, η επιλογή της τάσης Vc γίνεται με βάση το γεγονός πως η ακρίβεια της τάσης εισόδου Vin είναι πρακτικά πεπερασμένη ($\approx 0.5 mV$) οπότε είναι επιθυμητό η χλίση της σιγμοειδούς εξόδου να είναι όσο το δυνατόν μικρότερη



Σχήμα 3.1.4: Συνάρτηση μεταφοράς απλού κυκλώματος σιγμοειδούς για διαφορετικές τιμές Ibias

ώστε να αξιοποιείται βέλτιστα όλο το εύρος των δυνατών τάσεων εισόδου (-0.3V) με (0.3V). Έτσι η τάση Vc ορίζεται στα 100 mV ώστε παράλληλα να μην παρατηρείται σημαντική παραμόρφωση της εξόδου.

Αναφοριχά με την αναλυτική χυχλωματική ανάλυση του εν λόγω χυχλώματος σιγμοειδούς, θα ισχύουν τα αχόλουθα: Αρχικά, για την λειτουργία στην περιοχή υποκατωφλίου, το ρεύμα εκπομπού ενός NMOS τρανζίστορ, όπως αναφέρθηκε στην υποενότητα 2.2.3, θα δίνεται από την σχέση $I_D = I_0 \cdot e^{\frac{V_{GS}}{\xi V_T}}$. Όμως, για την ανάλυση της επίδρασης του φαινομένου σώματος στην λειτουργία του χυχλώματος της σιγμοειδούς είναι προτιμότερο να ξεκινήσω με τον παρακάτω πιο αναλυτικό τύπο για την λειτουργία στην περιοχή υποκατωφλίου [47]:

$$I_D = \frac{W}{L} \cdot I_t \cdot e^{\frac{V_{GS} - V_t}{\xi V_T}} \cdot (1 - e^{-\frac{V_{DS}}{V_T}})$$
(3.1.1)

Για $V_{DS} >> V_T$ ο τύπος αυτός απλοποιείται στην παραχάτω μορφή:

$$I_D = \frac{W}{L} \cdot I_t \cdot e^{\frac{V_{GS} - V_t}{\xi V_T}}$$

Όπου I_t είναι μια τιμή ρεύματος που εξαρτάται από την τεχνολογία κατασκευής και τις συνθήκες λειτουργίας. Ακόμη για απλότητα θεωρώ $I_c = \frac{W}{L} \cdot I_t \cdot e^{\frac{-V_t}{\xi V_T}}$ και έτσι



Σχήμα 3.1.5: Σχεδιάγραμμα κυκλώματος σιγμοειδούς συνάρτησης με διαφορικό ζεύγος διαφοράς

θα έχω:

$$I_D = I_c \cdot e^{\frac{V_{GS}}{\xi V_T}}$$

Έτσι για κάθε ένα από τα τρανζίστορ του διαφορικού ζεύγους διαφοράς θα ισχύει:

$$I_{D11} = I_{c11} \cdot e^{\frac{Vin - V_{S1}}{\xi V_T}}$$
$$I_{D12} = I_{c12} \cdot e^{\frac{-V_{S1}}{\xi V_T}}$$
$$I_{D13} = I_{c13} \cdot e^{\frac{Vin - V_{S2}}{\xi V_T}}$$
$$I_{D14} = I_{c14} \cdot e^{\frac{-V_{S2}}{\xi V_T}}$$



Σχήμα 3.1.6: Εξάρτηση εξόδου ρεύματος του κυκλώματος της εικόνα
ς3.1.5από την τάση $\rm Vc$

Όμως, όλα αυτά τα τρανζίστορ έχουν ίδιες διαστάσεις και τα ζευγάρια (M11, M13) και (M12, M14) έχουν ίδιες τάσεις V_B ενώ παράλληλα λόγω της λειτουργικής συμμετρίας του κυκλώματος (τα διοδικά συνδεδεμένα τρανζίστορ M21 και M22 έχουν ίδιες διαστάσεις)² οι τάσεις V_{S1} και V_{S2} θα είναι επίσης προσεγγιστικά ίσες. Συνεπώς και οι τάσεις V_t των παραπάνω ζευγαριών NMOS θα είναι ίσες. Άρα θα ισχύει $I_{c11} = I_{c14} = I_{c1}$ και $I_{c12} = I_{c13} = I_{c2}$ και οι παραπάνω σχέσεις των I_D απλοποιούνται ως εξής:

$$I_{D11} = I_{c1} \cdot e^{\frac{Vin - V_{S1}}{\xi V_T}}$$
$$I_{D12} = I_{c2} \cdot e^{\frac{-V_{S1}}{\xi V_T}}$$
$$I_{D13} = I_{c2} \cdot e^{\frac{Vin - V_{S1}}{\xi V_T}}$$
$$I_{D14} = I_{c1} \cdot e^{\frac{-V_{S1}}{\xi V_T}}$$

Αχόμη υφίσταται οι εξής δύο περιορισμοί:

$$I_{D11} + I_{D12} = I \tag{3.1.2}$$

 $^{^2{\}rm H}$ λειτουργική αυτή συμμετρία ισχύει απόλυτα όταν $V_c=-0.3V$ αλλά και για μεγαλύτερες τάσεις $V_c,$ δεδομένου της σχετικά ασθενούς επίδρασης του φαινομένου σώματος στην λειτουργία των NMOS, μπορούμε να κάνουμε την εύλογη υπόθεση ότι αυτή η συμμετρία εξακολουθεί να ισχύει σε μεγάλο βαθμό

$$I_{D13} + I_{D14} = I \tag{3.1.3}$$

Όπου $I(\propto Ibias)$ είναι το ρεύμα πόλωσης που τροφοδοτείται σε κάθε ένα από τα δύο διαφορικά ζεύγη από κάθε έναν από του όμοιους κασκοδικούς καθρέπτες ρεύματος (M01, M02, M03, M04) και (M01, M02, M03, M04). Παράλληλα διαιρώντας τα I_{D11}, I_{D12} και I_{D13}, I_{D14} αντίστοιχα θα έχω:

$$\frac{I_{D11}}{I_{D12}} = \frac{I_{c1}}{I_{c2}} \cdot e^{\frac{Vin}{\xi V_T}} = c_1 \cdot e^{\frac{Vin}{\xi V_T}}, \ c_1 = \frac{I_{c1}}{I_{c2}} = e^{\frac{V_{t12} - V_{t11}}{\xi V_T}}$$
(3.1.4)
$$\frac{I_{D13}}{I_{D14}} = \frac{I_{c2}}{I_{c1}} \cdot e^{\frac{Vin}{\xi V_T}} = \frac{1}{c_1} \cdot e^{\frac{Vin}{\xi V_T}}$$

Έτσι με βάση τις παραπάνω δύο σχέσεις και τις σχέσεις (3.1.2) και (3.1.3) προκύπτουν τα παρακάτω:

$$c_1 \cdot I_{D12} \cdot e^{\frac{Vin}{\xi V_T}} + I_{D12} = I \Longrightarrow I_{D12} = \frac{I}{1 + c_1 \cdot e^{\frac{Vin}{\xi V_T}}}$$
(3.1.5)

$$(1/c_1) \cdot I_{D14} \cdot e^{\frac{Vin}{\xi V_T}} + I_{D14} = I \Longrightarrow I_{D14} = \frac{I}{1 + (1/c_1) \cdot e^{\frac{Vin}{\xi V_T}}}$$
(3.1.6)

Η έξοδος αυτού του χυχλώματος σιμγοειδούς προχύπτει από τον χαθρεπτισμό του ρεύματος I_{D22} (με λόγο χαθρεπτισμού 1/2). Έτσι θα ισχύει

$$Iout = \frac{I_{D12} + I_{D14}}{2} = (I/2) \cdot \left(\frac{1}{1 + c_1 \cdot e^{\frac{Vin}{\xi V_T}}} + \frac{c_1}{c_1 + e^{\frac{Vin}{\xi V_T}}}\right) =>$$
$$Iout = (I/2) \cdot \frac{2c_1 + (c_1 + 1) \cdot e^{\frac{Vin}{\xi V_T}}}{(c_1 + e^{\frac{Vin}{\xi V_T}}) \cdot (1 + c_1 \cdot e^{\frac{Vin}{\xi V_T}})}$$
(3.1.7)

Αν θεωρήσουμε τώρα την ειδική περίπτωση όπου $V_c = -0.3V => V_{t11} = V_{t12} => c_1 = 1$, (από την σχέση (3.1.4)) η σχέση (3.1.7) απλοποιείται ως εξής:

$$Iout = \frac{I}{1 + e^{\frac{Vin}{\xi V_T}}}$$

η οποία σχέση περιγράφει μια ανεστραμμένη σιγμοειδή χαμπύλη.

Παρατηρώντας σε αυτό το σημείο την εξίσωση που περιγράφει το φαινόμενο σώματος (υποενότητα 2.2.2) διαφαίνεται πως όσο η τιμή της τάσης $Vc = V_{B12}$ αυξάνεται, θα τάση V_{t11} θα μειώνεται ενώ η τάση V_{t12} θα μένει σταθερή. Συνεπώς, από την εξίσωση (3.1.4) προχύπτει πως όσο αυξάνεται η τάση Vc θα αυξάνεται χαι η παράμετρος c_1 . Έτσι, μέσω των εξισώσεων (3.1.5), (3.1.6) που περιγράφουν τα ρεύματα των τρανζίστορ I_{D12} χαι I_{D14} αντίστοιχα χαι με την χρήση για το ξ

της ενδεικτικής τιμής 1.5, οι κανονικοποιημένες τιμές αυτών των ρευμάτων και της εξόδου (που προκύπτει ως το ημιάθροισμα τους) παρατίθενται στο Σχήμα 3.1.9.

Για την επιλογή διαστάσεων των τρανζίστορ χρειάζεται να ληφθούν υπόψη τα εξής: Αρχικά στόχος της σχεδίασης είναι το κύκλωμα, ως αναλογικός πολλαπλασιαστής, να παρουσιάζει όσο το δυνατόν υψηλότερη ακρίβεια και όσο το δυνατόν χαμηλότερη κατανάλωση ισχύος. Παράλληλα όμως με την κατανάλωση ισχύος, ένα εξίσου σημαντικό χαρακτηριστικό είναι η κατανάλωση ενέργειας ανά υπολογισμό. Αυτή, πέρα από την ισχύ εξαρτάται και από την μέγιστη συχνότητα λειτουργίας του κυκλώματος (με σχέση αντίστοφης αναλογίας) την οποία είναι επιθυμητό να μεγιστοποιήσουμε. Γενικά όμως για την επίτευξη υψηλής ακρίβειας (υψηλή γραμμικότητα και απουσία ταλαντώσεων) και υψηλής συχνότητας λειτουργίας απαιτείται σχετικά μεγάλο ρεύμα πόλωσης³. Επίσης για την πρώτη συνθήκη απαιτούνται σχετικά μεγάλες διαστάσεις τρανζίστορ και για την δεύτερη σχετικά μικρές. Παράλληλα, για να επιτευχθεί χαμηλή κατανάλωση ισχύος είναι απαραίτητη η χρήση χαμηλών ρευμάτων πόλωσης. Έτσι, με βάση αυτό το σύνολο συμβιβασμών οι προτεινόμενες διαστάσεις των τρανζίστορ για το εν λόγω κύκλωμα σιγμοειδούς παρουσιάζονται στον Πίνακα 3.1.

MOS transistors	$\mathbf{W/L} \; (\mu m/\mu m)$
M_{01}, M_{02}	0.4/1.6
M_{03}, M_{05}	0, 6/1.8
M_{04}, M_{06}	0,75/1.8
M_{11} - M_{14}	0.4/0.4
M_{21}, M_{22}	2.0/4.0
M_{23}	0.4/1.6

Πίναχας 3.1: Διαστάσεις τρανζίστορ (Σχήμα 3.1.5)

Εξετάζοντας την σχέση μεταξύ του Ιουτ και του Ibias για διαφορετικές τιμές Vin (Σχήμα 3.1.7), παρατηρώ πως η γραμμή που διαγράφεται, διέρχεται αρκετά πιο κοντά από το σημείο (0,0), όμως για πολύ μικρές τιμές Ibias εξακολουθεί να παρατηρείται μια μη γραμμική συμπεριφορά. Παράλληλα, το κύκλωμα αυτό μπορεί να χρησιμοποιηθεί ως αναλογικός πολλαπλασιαστής μόνο για θετικούς συντελεστές πολλαπλασιασμού. Για να υλοποιήσω και τα δύο πρόσημα χρησιμοποιώ ένα κύκλωμα (Σχήμα 3.1.10) το οποίο οδηγεί μια είσοδο ρεύματος από έναν διαφορετικό «δρόμο» για το εκάστοτε πρόσημο. Συγκεκριμένα, η επιλογή αυτή του επιθυμητού δρόμου σήματος εξόδου γίνεται με δύο αντίθετα σήματα τάσης ελέγχου (Vcnt1, Vcnt2 = -Vcnt1) τα οποία εφαρμόζονται στην θέση της τάσης αρνητικής τροφοδοσίας σε δύο κασκοδικούς NMOS καθρέπτες ρεύματος. Αν το σήμα που εφαρμόζεται είναι ίσο με την

³Το ρεύμα πόλωσης που τροφοδοτεί τα τρανζίστορ του διαφορικού ζεύγους διαφοράς επηρεάζεται και από τους λόγους καθρεπτισμού των σχετικών κασκοδικών καθρεπτών ρεύματος.

αρνητική τροφοδοσία, ο καθρέπτης θα λειτουργεί κανονικά, ενώ για τάση ελέγχου ίση με την θετική τροφοδοσία αυτός θα βρίσκεται σε αποκοπή.

Για την επεξήγηση της λειτουργίας του εχάστοτε χαθρέπτη ρεύματος ως ελεγχόμενης πύλης μετάδοσης είναι χρήσιμο να εξετάσουμε τον τύπο (1) που περιγράφει αναλυτιχά την λειτουργία ενός NMOS τρανζίστορ στην περιοχή υποχατωφλίου (ο τύπος για την περίπτωση του PMOS έχει απλώς V_{SG} στην θέση της τάσης V_{GS}). Συγχεχριμένα για την περίπτωση που στην θέση του εχάστοτε σήματος ελέγχου εφαρμόζεται η τάση Vss, θα ισχύει για όλα τα τρανζίστορ του χαθρέπτη $V_{DS} >> V_T$ χαι ο τύπος (1) θα απλοποιείται στην μορφή:

$$I_D = \frac{W}{L} \cdot I_t \cdot e^{\frac{V_{GS} - V_t}{\xi V_T}}$$

Έτσι, για τα δύο τρανζίστορ του NMOS χαθρέπτη που έχουν τους αχροδέχτες τους στην αρνητική τάση τροφοδοσίας, εφόσον και οι τάσεις V_G είναι κοινές, το ρεύμα του δεξιού μέλους του καθρέπτη θα είναι ίσο με αυτό του αριστερού μέρους επί τον λόγο καθρεπτισμού $m = \frac{(W_2/L_2)}{(W_1/W_1)}$ (ο οποίος σε αυτή την περίπτωση είναι ίσος με 1), όπου η ένδειξη 2 υποδηλώνει το συγκεκριμένο τρανζίστορ του δεξιού μέρους και η ένδειξη 1 υποδηλώνει το αριστερό τρανζίστορ. Σε περίπτωση όμως που εφαρμοστεί το σήμα ελέγχου 0.3V η συνθήκη $V_{DS} >> V_T$ προφανώς παύει να ισχύει για τα τρανζίστορ του δεξιού μέρους και αντί αυτού οι αντίστοιχες τάσεις V_{DS} θα τείνουν στο 0. Άρα τα αντίστοιχα τρανζίστορ θα μπαίνουν σε αποχοπή και το σήμα εισόδου δεν θα μεταδίδεται μέσω αυτών.

Έτσι η λειτουργία αυτού του χυχλώματος προσήμου περιγράφεται από την αχόλουθη εξίσωση:

$$Iout = \begin{cases} Iin, & Vcnt1 = -0.3V \\ -Iin, & Vcnt1 = 0.3V \end{cases}$$

Αχόμη, για την ορθή λειτουργία αυτού του χυχλώματος, δεδομένου ότι στην έξοδο του καταλήγει PMOS και NMOS καθρέπτης ρεύματος, η τάση του κόμβου εξόδου θα πρέπει να βρίσκεται όσο το δυνατόν πιο κοντά στο μέσο της τροφοδοσίας (0V σε αυτήν την περίπτωση).

Πίναχας 3.2: Διαστάσεις τρανζίστορ (Σχήμα 3.1.10)

MOS transistors	$\mathbf{W/L}~(\mu m/\mu m)$
M_{31}, M_{32}	0.4/0.8
M_{33} - M_{36}	0.4/0.2
M_{41} - M_{44}	0.4/3.2



Σχήμα 3.1.7: Παράδειγμα εξόδου κυκλώματος του Σ
χήματος 3.1.5 συναρτήσει του Ibias για ${\rm Vr}=0$

3.2 Κυχλώματα πολλαπλασιαστών

Στην προηγούμενη ενότητα παρουσιάστηκαν δύο κυκλώματα σιγμοειδούς συνάρτησης και εξετάστηκε η χρήση τους ως αναλογικούς πολλαπλασιαστές. Όπως φάνηκε όμως από τις σχετικές προσομοιώσεις, η υπολογιστική ακρίβεια που αυτά μπορούσαν να επιτύχουν ήταν αρκετά περιορισμένη. Έτσι, στην παρούσα ενότητα εξετάζεται ο σχεδιασμός πρακτικών κυκλωμάτων αναλογικών πολλαπλασιαστών βασιζόμενων στο κύκλωμα του Σχήματος 3.1.5, με σημαντικά βελτιωμένη συμπεριφορά.

3.2.1 Κύκλωμα πολλαπλασιαστή συνελικτικών φίλτρων

Από την ανάλυση της προηγούμενης υποενότητας χατέστη εμφανές πως το δεύτερο χύχλωμα σιγμοειδούς (με διαφοριχό ζεύγος διαφοράς) που παρουσιάστηχε είχε σημαντιχά χαλύτερα χαραχτηριστιχά από την πρώτη, απλούστερη υλοποίηση. Όμως αυτό εξαχολουθούσε να παρουσιάζει αισθητή μη γραμμιχότητα για πολύ χαμηλές τιμές Ibias. Με βάση λοιπόν αυτήν την παρατήρηση, ένας τρόπος για να υλοποιηθεί ένας αναλογιχός πολλαπλασιαστής με βελτιωμένη γραμμιχότητα είναι να αποφύγουμε την λειτουργία αυτού του χυχλώματος σιγμοειδούς για πολύ χαμηλές τιμές ρεύματος εισόδου. Ωστόσο, ένα πραχτιχό χύχλωμα αναλογιχού πολλαπλασιαστή θα πρέπει να μπορεί να λειτουργήσει για τιμές εισόδου χοντά στο μηδέν το οποίο συνεπάγεται



Σχήμα 3.1.8: Συνάρτηση μεταφοράς χυχλώματος του Σχήματος 3.1.5 για διαφορετιχές τιμές Ibias

πολύ χαμηλές τιμές αντίστοιχων ρευμάτων . Έτσι, για να λυθεί αυτό το πρόβλημα, προτείνεται η υλοποίηση του Σχήματος 3.2.1.

Πίναχας 3.3: Διαστάσεις επιπρόσθετων τρανζίστορ Σχήματος 3.2.1

MOS transistors	$\mathbf{W/L}~(\mu m/\mu m)$
M_{51}, M_{52}	0.4/1.6
M_{61} - M_{64}	0.4/1.6
M_{71}, M_{72}	0.4/1.6

Συγκεκριμένα, η υλοποίηση αυτή περιλαμβάνει δύο πανομοιότυπα (με ίδιες διαστάσεις τρανζίστορ) κυκλώματα σιγμοειδούς με ελεγχόμενο πρόσημο (μπλε ορθογώνια), τα οποία λειτουργούν ως επιμέρους πολλαπλασιαστές για να προκύψει τελικά το επιθυμητό γινόμενο. Αρχικά στο επάνω κύκλωμα σιγμοειδούς τροφοδοτείται ως είσοδος ένα μικρό ρεύμα πόλωσης (*Ib_lin*) του οποίου η τιμή επιλέγεται έτσι ώστε αυτό το κύκλωμα να λειτουργεί οριακά εντός της επιθυμητής γραμμικής περιοχής. Παράλληλα όμως μέσω των καθρεπτών ρεύματος (M01, M02, M51, M52) και (M61, M62, M63, M64), οι οποίοι έχουν λόγο καθρεπτισμού 1, το ρεύμα αυτό τροφοδοτείται στην είσοδο του καθρέπτη (M71, M72). Εκεί αυτό αθροίζεται με το ρεύμα εισόδου του συνολικού κυκλώματος και το προκύπτον άθροισμα ρεύματος τροφο Σχήμα 3.1.9: Κανονικοποιημένες τιμές των $I_{D12}, I_{D14}, Iout$ για διαφορετικές τιμές V_c όπως προκύπτουν απο τις αναλυτικές σχέσεις (3.1.4) και (3.1.5)



δοτείται ως είσοδος στο δεύτερο κύκλωμα σιγμοειδούς. Τα δύο αυτά κυκλώματα σιγμοειδούς έχουν αντίθετα σήματα τάσεων ελέγχου (Vcnt1, Vcnt2) και συνεπώς υλοποιούν αντίθετα πρόσημα. Ακολούθως, η έξοδος του συνολικού κυκλώματος προκύπτει από τον κοινό κόμβο εξόδου των δύο κυκλωμάτων σιγμοειδών όπου οι έξοδοι ρεύματος τους αφαιρούνται. Έτσι, για τον πολλαπλασιασμό ενός θετικού ρεύματος εισόδου Iin με κάποιο συντελεστή α, υπολογίζονται τα γινόμενα $\alpha \cdot Ib_lin$ και $\alpha \cdot (Ib_lin + Iin)$ και στη συνέχεια αυτά αφαιρούνται στον χόμβο εξόδου για να παραχθεί το τελικό επιθυμητό γινόμενο $\alpha \cdot Iin$. Με αυτόν τον τρόπο εξασφαλίζω ότι και τα δύο κυκλώματα σιγμοειδών θα λειτουργών εντός της γραμμικής περιοχής για όλες τις δυνατές (θετικές) τιμές ρευμάτων εισόδου. Από τις αντίστοιχες προσομοιώσεις (Σχήματα 3.2.2 και 3.2.3) επαληθεύεται και η σημαντική βελτίωση της γραμμικότητας, συνεπώς και της υπολογιστικής ακρίβειας, που αυτό το χύκλωμα προσφέρει.

Εκτός από την σχεδίαση του κυκλώματος, για τη χρήση του σε πρακτικές εφαρμογές απαιτείται και η αντιστοίχιση ενός δεδομένου πραγματικού συντελεστή πολλαπλασιασμού με ένα διάνυσμα εισόδου (Vin, Vent1, Vent2). Αρχικά, όπως διαφαίνεται και από τις αντίστοιχες προσομοιώσεις (Σχήμα 3.2.3), ο μέγιστος (κατά απόλυτη τιμή) συντελεστής πολλαπλασιασμού που μπορεί να υλοποιήσει το κύκλωμα είναι περίπου 2.1. Όμως, δεδομένου ότι το εν λόγω κύκλωμα μπορεί να υλοποιήσει συντελεστές πολύ κοντά στο 0, με την προσθήκη ενός κυκλώματος κλιμάκωσης ε-



Σχήμα 3.1.10: Σχεδιάγραμμα κυκλώματος προσήμου

ξόδου (υποενότητα 3.2.2) αυτό μπορεί να χρησιμοποιηθεί για ένα πολύ μεγαλύτερο εύρος συντελεστών. Αναφορικά με την τάση Vcnt1, αυτή θεωρείται 0.3V για την περίπτωση θετικού προσήμου και -0.3V αλλιώς (και αντίθετα για την Vcnt2). Για την αντιστοίχιση ενός συντελεστή με την τάση Vin, η οποία ρυθμίζει το πλάτος εξόδου, είναι απαραίτητη η χρήση μίας συγκεκριμένης συνάρτησης μεταφοράς (ανεξάρτητης του ρεύματος εισόδου), καθώς το πλάτος του *Iin* μπορεί να πάρει μεγάλο εύρος τιμών. Η συνάρτηση αυτή λαμβάνεται από τον σταθμισμένο μέσο όρο των συναρτήσεων μεταφοράς για ρεύμα εισόδου ίσο με 2nA, 4nA, 6nA, 8nA και 10nA και είναι διαφορετική για το θετικό και το αρνητικό πρόσημο, καθώς παρατηρούνται μικρές αποκλίσεις μεταξύ τους. Μια από τις βασικότερες τέτοιες αποκλίσεις είναι το γεγονός ότι η συνάρτηση μεταφοράς παίρνει ελαφρώς μεγαλύτερες τιμές (κατά πλάτος) στην περίπτωση των θετικών συντελεστών. Εδώ να σημειωθεί πως αυτή η διαφορά θα μπορούσε θεωρητικά να διορθωθεί με μια μικρή προσαρμογή των διαστάσεων των τρανζίστορ. Όμως η απαιτούμενη προσαρμογή θα ήταν στην ίδια τάξη μεγέθους



Σχήμα 3.1.11: Σχεδιάγραμμα κυκλώματος σιγμοειδούς ακολουθούμενο από το κύκλωμα προσήμου

με τις τυχαίες αποκλίσεις PVT και συνεπώς για την αντιμετώπιση του σφάλματος που εισάγεται στο τελικό αποτέλεσμα χρησιμοποιώ μια πιο ευέλικτη προσέγγιση που περιγράφεται στην υποενότητα 5.2.2. Επιπλέον, το κύκλωμα αυτό επιτυγχάνει εξαιρετικά χαμηλή κατανάλωση ισχύος, καθώς μπορεί να χρησιμοποιηθεί για εισόδους της τάξης των λίγων εκατοντάδων pA ενώ μπορεί να λειτουργήσει σε συχνότητες μέχρι περίπου τα 250kHz χωρίς αισθητή απώλεια υπολογιστικής ακρίβειας.

3.2.2 Κύκλωμα κλιμάκωσης ρεύματος εξόδου

Το χύχλωμα χλιμάχωσης εξόδου (Σχήμα 3.2.4) είναι ουσιαστιχά μια απλουστευμένη έχδοση του αναλογιχού πολλαπλασιαστή που παρουσιάστηχε λεπτομερώς παραπάνω και αποσχοπεί στην χλιμάχωση της εξόδου ρεύματος ενός προηγούμενου συστήματος χατά μία σταθερή τιμή, η οποία ορίζεται ανάλογα με την επιθυμητή εφαρμογή. Συγχεχριμένα, αυτό χρησιμοποιεί δύο χυχλώματα σιγμοειδούς, όπου σε χάθε ένα διατηρείται ο κλάδος μόνο του ενός προσήμου (του θετικού για το επάνω κύκλωμα και του αρνητικού για το κάτω), εξαλείφοντας παράλληλα την ανάγκη για τα σήματα ελέγχου Vcnt1, Vcnt2. Ακόμη, όπως φαίνεται και στο Σχήμα 3.2.4, η είσοδος ρεύματος είναι τοποθετημένη λίγο διαφορετικά ώστε να δέχεται ρεύμα εισόδου αντίθετης φοράς από τον πολλαπλασιαστή της προηγούμενης υποενότητας. Επίσης στους ακροδέκτες της πηγής και του σώματος των τρανζίστορ M61 και M63 εφαρμόζεται σήμα τάσης Vdm (Vdm > Vdd) αντί για την θετική τάση τροφοδοσίας, ώστε να διευρυνθεί το εύρος των τάσεων του κόμβου εισόδου για το οποίο λειτουργεί ο συγκεκριμένος καθρέπτης ρεύματος. Σε αντίθετη περίπτωση, από τις προσομοιώσεις παρατηρείται μη επιθυμητή συμπίεση της εξόδου ρεύματος, η οποία εισάγει σημαντικό σφάλμα. Αναφορικά με τις διαστάσεις των τρανζίστορ του κυκλώματος κλιμάκωσης εξόδου, αυτές είναι ίδιες με το κύκλωμα του πολλαπλασιαστή αναλογικών φίλτρων για όλα τα τρανζίστορ εκτός των M61 – M64 που έχουν όλα διαστάσεις $W/L = 0.4 \mu m/3.2 \mu m$.

3.3 Κύκλωμα ReLU συνάρτησης ενεργοποίησης

Για την χυχλωματιχή υλοποίηση της συνάρτησης ενεργοποίησης ReLU είναι επιθυμητή η ανεμπόδιστη διέλευση των ρευμάτων θετιχής φοράς (εισερχόμενων) και η αποχοπή των ρευμάτων αρνητιχής φοράς (εξερχόμενων). Αυτές οι επιθυμητές προδιαγραφές (*Iout = max(Iin, 0*)) περιγράφουν την λειτουργία του ιδανιχού NMOS καθρέπτη ρεύματος. Για την πραγματιχή σχεδίαση λοιπόν επιλέγεται ένας χασχοδιχός NMOS χαθρέπτης ρεύματος εξαιτίας του πιο σταθερού λόγου χαθρεπτισμού που προσφέρει. Σε περίπτωση όμως που το χύχλωμα αυτό συνδέεται κατευθείαν στην έξοδο ενός συνελιχτιχού φίλτρου (όπως συμβαίνει στο χύχλωμα της ενότητας 4.4) πρέπει να διασφαλιστεί ότι το χύχλωμα της ReLU λειτουργεί για ένα μεγάλο εύρος τάσεων του χόμβου εισόδου ώστε να μην παραμορφώνεται η έξοδος ρεύματος του φίλτρου. Στην παρούσα εφαρμογή, αυτό επιτυγχάνεται με την χρήση ενός ξεχωριστού σήματος τάσης Vsm για τους αχροδέχτες της πηγής των δύο χάτω τρανζίστορ και του αχροδέχτες σώματος όλων των τρανζίστορ του χαθρέπτη. Για την περίπτωση που το χύχλωμα της ReLU χρησιμοποιείται ξεχωριστά η τάση Vsm ορίζεται ίση με την αρνητιχή τάση τροφοδοσίας.

Πίναχας 3.4: Διαστάσεια	ς τρανζίστορ	Σχήματος	3.3.1
-------------------------	--------------	----------	-------

MOS transistors	$\mathbf{W/L} \; (\mu m/\mu m)$
M_{r1}, M_{r2}	0.8/1.2
M_{r3}, M_{r4}	0.4/4.0

3.4 Κυχλώματα Winner-Take-All

Για την παραγωγή της εξόδου ενός μοντέλου μηχανικής μάθησης που εκτελεί κάποια διαδικασία ταξινόμησης απαιτείται ενα σύστημα που θα μπορεί να ελέγχει για ποία από τις εκάστοτε κλάσεις αυτό παρουσιάζει εντονότερη ενεργοποίηση. Αυτό συνήθως ισοδυναμεί με την εύρεση της εξόδου του μοντέλου που έχει την μεγαλύτερη τιμή, δηλαδή με τον μαθηματικό τελεστή argmax. Στην παρούσα ενότητα εξετάζονται δύο σχετικές κυκλωματικές υλοποιήσεις και επεξηγείται το σκεπτικό της τελικής επιλογής μεταξύ αυτών και άλλων υλοποιήσεων.

3.4.1 Απλό Winner-Take-All χύχλωμα

Σε χυχλωματικό επίπεδο μια απλή και ευρέως χρησιμοποιούμενη υλοποίηση αυτού του τελεστή είναι το χύχλωμα Winner-Take-All (WTA) που προτάθηκε από τον J. Lazzaro [49] (Σχήμα 3.4.1). Το χύχλωμα αυτό δέχεται N εισόδους ρεύματος και ιδανικά η έξοδος του χλάδου με το μέγιστο ρεύμα εισόδου θα ισούται με το ρεύμα πόλωσης, ενώ όλες οι άλλες έξοδοι ρεύματος θα είναι μηδενικές. Για να εξηγηθεί σε μεγαλύτερη λεπτομέρεια η λειτουργία αυτού του χυχλώματος, θα εξεταστεί η περίπτωση των δύο εισόδων (Σχήμα 3.4.2). Συγχεκριμένα, όσο τα τρανζίστορ M11 και M21 λειτουργούν στην περιοχή υποκατωφλίου, το ρεύμα τους θα περιγράφεται από τις παρακάτω εξισώσεις (το I_{c1} είναι κοινό γιατί αυτά τα τρανζίστορ έχουν ίδιες διαστάσεις):

$$I_{D11} = I_{c1} \cdot e^{\frac{V_1 - V_{ss}}{\xi V_T}} \cdot (1 - e^{-\frac{V_{D11} - V_{ss}}{V_T}})$$
$$I_{D21} = I_{c1} \cdot e^{\frac{V_1 - V_{ss}}{\xi V_T}} \cdot (1 - e^{-\frac{V_{D21} - V_{ss}}{V_T}})$$

όπου V₁ είναι τη τάση του κοινού κόμβου πύλης.

Στην περίπτωση που τα δύο ρεύματα εισόδου είναι ίδια, τότε λόγω συμμετρίας ⁴ χαι οι δύο έξοδοι ρεύματος θα είναι ίσες με το μισό του ρεύματος πόλωσης (Ibias_wta). Αν τώρα θεωρήσουμε πως, ξεκινώντας από την χατάσταση ηρεμίας ($I_{in1} = I_{in2}$) χαι χωρίς βλάβη της γενιχότητας, το I_{in1} αυξηθεί ($I_{in1} = I_{in2} + \delta$), θα αυξηθεί χαι η τάση πύλης του M11 (V_1). Η τάση αυτή της πύλης είναι χοινή με το τρανζίστορ M21 συνεπώς θα αυξηθεί και η τάση V_{GS21} που από τον προηγούμενο τύπο για το I_{D21} διαφαίνεται πως προχαλεί μια αυξητική τάση σε αυτό το ρεύμα. Έτσι δεδομένου πως το ρεύμα I_{D21} παραμένει σταθερό και η συνάρτηση αυτού του ρεύματος είναι γνησίως αύξουσα ως προς την τάση V_{D21} ($\frac{\partial I_{D21}}{\partial V_{D21}} = \frac{V_{D21} - Vss}{V_T} \cdot I_{c1} \cdot e^{\frac{V_1 - Vss}{\xi V_T}} \cdot e^{-\frac{V_{D21} - Vss}{V_T}} > 0$), για μικρές τιμές δ, το M21 μπορεί να προσαρμοστεί μειώνοντας την τάση V_{D21} . Σε αυτήν την περίπτωση η

⁴Εδώ θεωρούμε πως και οι δύο κλάδοι (γενικά Ν) είναι πανομοιότυποι και συνεπώς τα τρανζίστορ τους στις ίδιες θέσεις έχουν τις ίδιες διαστάσεις.

μείωση της $V_{D21} = V_{G22}$ θα προκαλέσει και μείωση του ρεύματος I_{D22} που περιγράφεται από τη σχέση: $I_{D22} = I_{c2} \cdot e^{\frac{V_{G22}-V_1}{\xi V_T}}$ το οποίο συνεπάγεται και αύξηση του ρεύματος I_{D12} δεδομένου ότι $I_{D12} + I_{D22} = Ibias_wta$ Συγκεκριμένα, για την περίπτωση αυτή των μικρών δ θα ισχύουν αναλυτικά τα ακόλουθα:

$$\frac{I_{D11}}{I_{D21}} = \frac{\left(1 - e^{-\frac{V_{D11} - V_{ss}}{V_T}}\right)}{\left(1 - e^{-\frac{V_{D21} - V_{ss}}{V_T}}\right)} =>$$

$$1 - e^{-\frac{V_{D21} - V_{ss}}{V_T}} = c_1 \cdot \left(1 - e^{-\frac{V_{D11} - V_{ss}}{V_T}}\right) =>$$

$$V_{D21} = V_{ss} - V_T \cdot ln(1 - c_1 \cdot \left(1 - e^{-\frac{V_{D11} - V_{ss}}{V_T}}\right))$$

όπου $c_1 = \frac{I_{D21}}{I_{D11}}$

Παράλληλα για τα τρανζίστορ M12, M22 ισχύουν οι προσεγγιστικές σχέσεις:

$$I_{D12} = I_{c2} \cdot e^{\frac{V_{D11} - V_1}{\xi V_T}}$$
$$I_{D22} = I_{c2} \cdot e^{\frac{V_{D21} - V_1}{\xi V_T}}$$

Από τις οποίες προχύπτει:

$$\frac{I_{D12}}{I_{D22}} = e^{\frac{V_{D11} - V_{D21}}{\xi V_T}}$$

Το οποίο σε συνδυασμό με τον περιορισμό $I_{D12} + I_{D22} = Ibias_wta$ συνεπάγεται ότι θα ισχύει:

$$I_{D22} = \frac{Ibias_wta}{1 + e^{\frac{V_{D11} - V_{D21}}{V_T}}}$$
$$I_{D12} = \frac{Ibias_wta \cdot e^{\frac{V_{D11} - V_{D21}}{V_T}}}{1 + e^{\frac{V_{D11} - V_{D21}}{V_T}}}$$

Για μεγαλύτερες όμως τιμές δ, η απαραίτητη πτώση της V_{D21} είναι τέτοια που θα αναγκάσει το M21 να βγει από την περιοχή υποκατωφλίου. Έτσι θα ισχύει $V_{D21} \rightarrow Vss => V_{G22} \rightarrow Vss$ και ως συνέπεια το M22 θα μπει σε αποκοπή. Άρα σε αυτήν την περίπτωση πρακτικά όλο το ρεύμα πόλωσης Ibias_wta θα διέρχεται από το τρανζίστορ M21 (Iout1 = Ibias_wta, Iout2 = 0) Έτσι η έξοδος αυτού του κυκλώματος θα έχει την συμπεριφορά που φαίνεται στο Σχήμα 3.4.3.

3.4.2 Κασχοδιχό Winner-Take-All χύχλωμα

Η διαχωριστική ικανότητα του προαναφερθέντος κυκλώματος μπορεί να βελτιωθεί περεταίρω με την προσθήκη ενός ακόμη συμπληρωματικού PMOS WTA κυκλώματος, όπως φαίνεται στο Σχήμα 3.4.4. Συγκεκριμένα το κύκλωμα αυτό δέχεται ως εισόδους τα ρεύματα εξόδου του απλού NMOS WTA κυκλώματος και έτσι συγκρίνει ρεύματα με μεγαλύτερες διαφορές από τις αρχικές εισόδους, γεγονός που συνεπάγεται την μείωση του εύρους της περιοχής όπου παρατηρούνται «πολλαπλοί νικητές» (Σχήμα 3.4.3). Εδώ να σημειώσω πως θεωρητικά θα μπορούσα να χρησιμοποιήσω ακόμη περισσότερα συμπληρωματικά WTA κυκλώματα, όμως από προσομοιώσεις παρατηρώ ότι με την προσθήκη του δεύτερου πρακτικά επιτυγχάνεται η βέλτιστη επίδοση. Επίσης, οι διαστάσεις όλων των τρανζίστορ του συγκεκριμένου κασκοδικού κυκλώματος WTA είναι ίσες με $W/L = 0.4/0.2(\mu m/\mu m)$ και τα δύο ρεύματα από κον παραχάτω τύπο:

$$Iout_{-i} = \begin{cases} 10nA, & if \ argmax(Iin_{-j}, \ \forall j \in \{1, N\}) = i \\ 0, & otherwise \end{cases}$$

Σχήμα 3.2.1: Σχεδιάγραμμα κυκλώματος πολλαπλασιαστή συνελικτικών φίλτρων (Κάθε ένα από τα μπλέ ορθογώνια χρησιμοποιείται στη θέση του κυκλώματος του Σχήματος 3.1.11 που βρίσκεται εντός του μωβ περιγράμματος)



Σχήμα 3.2.2: Παράδειγμα της απόλυτης εξόδου του χυχλώματος του Σχήματος 3.2.1 συναρτήσει του Iin για Vr = 0 και για τα δύο πρόσημα (το διάγραμμα εστιάζει στις χαμηλές τιμές Iin για να φανεί η σημαντικά βελτιωμένη γραμμικότητα αυτού του χυχλώματος)









Σχήμα 3.2.4: Σχεδιάγραμμα χυχλώματος κλιμάχωσης εξόδου ρεύματος





Σχήμα 3.4.1: Σχεδιάγραμμα Lazzaro WTA κυκλώματος Ν εισόδων





Σχήμα 3.4.2: Σχεδιάγραμμα Lazzaro WTA κυκλώματος 2 εισόδων

Σχήμα 3.4.3: Έξοδοι ρεύματος απλού και κασκοδικού WTA κυκλώματος 2
 εισόδων για Iin1=5nAκαι $Ibias_wta=10nA$





Σχήμα 3.4.4: Σχεδιάγραμμα κασκοδικού WTA κυκλώματος Ν εισόδων

Κεφάλαιο 4

Προτεινόμενη υλοποίηση αναλογικού νευρωνικού δικτύου

Στις προηγούμενες ενότητες εισήχθησαν οι απαραίτητες θεωρητικές έννοιες της μηχανικής μάθηση και ιδιαίτερα των νευρωνικών δικτύων καθώς επίσης παρουσιάστηκαν και τα βασικά δομικά κυκλώματα που απαιτεί η προτεινόμενη υλοποίηση. Στο παρόν κεφάλαιο αναλύεται ο τρόπος με τον οποίο όλα αυτά τα δομικά κυκλώματα συνδυάζονται για να προκύψει τελικά ένα ολοκληρωμένο σύστημα με την δυνατότητα αναλογικής υλοποίησης ποικίλων αρχιτεκτονικών συνελικτικών νευρωνικών δικτύων. Συγκεκριμένα, αρχικά παρουσιάζεται η αφηρημένη αρχιτεκτονική του προς υλοποίηση CNN και ακολούθως αναλύεται η κυκλωματική υλοποίηση κάθε ενός από τα επίπεδα του.

4.1 Αρχιτεκτονική προτεινόμενου συνελικτικού νευρωνικού δικτύου

Η προτεινόμενη αρχιτεκτονική συνελικτικού νευρωνικού δικτύου παρουσιάζεται στο Σχήμα 4.1.1. Συγκεκριμένα, το δίκτυο αυτό είναι αρχικά σχεδιασμένο για να λαμβάνει ως είσοδο 28 × 28 μονοχρωματικές εικόνες τις οποίες τροφοδοτεί στο πρώτο συνελικτικό του επίπεδο που περιλαμβάνει τρία 3 × 3 συνελικτικά φίλτρα. Οι τιμές των pixel αυτών των εικόνων μπορούν να προέρχονται απευθείας από έναν αισθητήρα ή από μια αναλογική μνήμη και σε κάθε περίπτωση τα αντίστοιχα ρεύματα θα παίρνουν θετικές τιμές. Εδώ να σημειώσω πως εικόνες τριών καναλιών μπορούν να μετατραπούν κυκλωματικά σε μονοχρωματικές με πολύ απλό τρόπο μέσω ενός κοινού κόμβου εξόδου και ενός καθρέπτη ρεύματος με κατάλληλο λόγο καθρεπτισμού. Στις εξόδους των τριών αυτών φίλτρων εφαρμόζεται η συνάρτηση ενεργοποίησης ReLU και στην συνέχεια ακολουθεί ένα επίπεδο average pooling που βελτιώνει την μεταθετική αμεταβλητότητα του δικτύου και επιτρέπει την χρήση ανώτερων στρωμάτων με μικρότερες διαστάσεις. Ο λόγος που σε αυτό το σημείο χρησιμοποιείται ως επίπεδο ομαδοποίησης ένα στρώμα average pooling αντί για max pooling, όπου το δεύτερο προσφέρει ελαφρώς καλύτερες επιδόσεις σε προβλήματα ταξινόμησης [50], είναι γιατί η υλοποίηση του πρώτου σε χυχλωματικό επίπεδο είναι πολύ απλούστερη. Συνεπώς η υλοποίηση του επιπέδου average pooling θα εισάγει μικρότερο σφάλμα, δεδομένου ότι για αυτήν απαιτείται απλώς ένα κύκλωμα κλιμάκωσης εξόδου. Ακολούθως το δεύτερο και το τρίτο επίπεδο του δικτύου αποτελούνται από τρία 3×3 συνελικτικά φίλτρα τριών καναλιών ενώ το τέταρτο επίπεδο αποτελείται από ένα τέτοιο φίλτρο, όπου σε όλα εφαρμόζεται πάλι ReLU ενεργοποίηση. Σε όλα τα προηγούμενα συνελικτικά επίπεδα δεν εφαρμόζεται padding και έτσι ξεκινώντας με μια 28×28 εικόνα, η τελική έξοδος αυτών έχει διαστάσεις
 $7\times 7.$ Ακολούθως, για την χεφαλή ταξινόμησης του διχτύου χρησιμοποιείται ένα πλήρως συνδεδεμένο επίπεδο με $7\times 7=49$ εισόδους και 43 εξόδους. Κατά την διαδικασία της εκπαίδευσης του διχτύου το πλήρως συνδεδεμένο επίπεδο αχολουθείται από ένα επίπεδο Softmax ενεργοποίησης και έπειτα η τελική έξοδος παράγεται μέσω του τελεστή argmax ενώ μετά την εκπαίδευση η Softmax ενεργοποίηση παραλείπεται. Αυτή η επιλογή γίνεται επειδή η εφαρμογή της συνάρτησης ενεργοποίησης Softmax δεν μεταβάλλει το αποτέλεσμα του τελεστή argmax ενώ σε κυκλωματικό επίπεδο προσθέτει σημαντική πολυπλοκότητα.

4.2 Κύκλωμα συνελικτικού φίλτρου 3×3 ενός καναλιού εισόδου με ReLU ενεργοποίηση

Ένα 3 × 3 συνελκτικό φίλτρο ενός καναλιού πρακτικά περιλαμβάνει 3 × 3 = 9 υπολογιστικές μονάδες πολλαπλασιασμού των οποίων η έξοδος αθροίζεται μαζί με μια τιμή bias. Έτσι αυτό μπορεί να υλοποιηθεί κυκλωματικά με την χρήση 9 αναλογικών πολλαπλασιαστών (όπως αυτοί περιγράφηκαν στην υποενότητα 3.2.1) και μίας ηλεκτρονικά ελεγχόμενης πηγής ρεύματος για την υλοποίηση του bias. Όλα αυτά τα κυκλώματα θα οδηγούν σε έναν κοινό κόμβο εξόδου όπου τα ρεύματα τους θα αθροίζονται για να προκύψει το αποτέλεσμα του φίλτρου. Ακολούθως, αυτή η έξοδος ρεύματος θα οδηγείται στο κύκλωμα υλοποίησης της συνάρτησης ενεργοποίησης ReLU, όπως παρουσιάζεται στο Σχήμα 4.1.2.

Για τον υπολογισμό της εξόδου του παραπάνω φίλτρου, αυτό πραχτικά «σύρεται» πάνω από την εικόνα εισόδου με μοναδιαίο βήμα (Σχήμα 4.2.2) και τα αποτελέσματα μπορούν να αποθηκευτούν για μετέπειτα χρήση σε κύτταρα αναλογικής μνήμης. Θεωρητικά για τον υπολογισμό του αποτελέσματος του φιλτραρίσματος με ένα δεδομένο συνελικτικό φίλτρο χρειάζεται μόνο ένα τέτοιο κύκλωμα, όμως αυτή η διαδικασία μπορεί να παραλληλοποιηθεί χρησιμοποιώντας πολλαπλά ίδια κυκλώματα που θα υπολογίζουν την δισδιάστατη συνέλιξη σε διαφορετικές περιοχές της εικόνας.



Σχήμα 4.1.1: Αρχιτεκτονική προτεινόμενου συνελικτικού νευρωνικού δικτύου

Επίσης, δεδομένου ότι οι συντελεστές πολλαπλασιασμού αυτών των συνελικτικών φίλτρων είναι ηλεκτρονικά ελεγχόμενοι, το ίδιο κύκλωμα μπορεί να χρησιμοποιηθεί για την υλοποίηση πολλαπλών διαφορετικών φίλτρων με ίδιες διαστάσεις.

4.3 Κύκλωμα συνελικτικού φίλτρου 3×3 τριών καναλιών εισόδου με ReLU ενεργοποίηση

Ένα 3×3 συνελικτικό φίλτρο τριών καναλιών πρακτικά περιλαμβάνει $3 \times 3 \times 3 = 27$ υπολογιστικές μονάδες πολλαπλασιασμού των οποίων η έξοδος θα αθροίζεται μαζί με μια τιμή bias. Έτσι, όμοια με το αντίστοιχο φίλτρο ενός καναλιού, αυτό μπορεί να υλοποιηθεί κυκλωματικά με την χρήση 27 αναλογικών πολλαπλασιαστών της υποενότητας 3.2.1 και μιας ελεγχόμενης πηγή ρεύματος. Ακολούθως, οι έξοδοι



Σχήμα 4.2.1: Σχεδιάγραμμα κυκλώματος 3 × 3 συνελικτικού φίλτρου ενός καναλιού με ReLU ενεργοποίηση

ρεύματος όλων αυτών των κυκλωμάτων αθροίζονται σε έναν κοινό κόμβο εξόδου, του οποίου η τάση διατηρείται κοντά στο 0, και η έξοδος ρεύματος που προκύπτει αποθηκεύεται σε αναλογική μνήμη. Στην συνέχεια οι τιμές αυτές ανακτώνται από την μνήμη και οδηγούνται στο κύκλωμα της συνάρτησης ενεργοποίησης ReLU από όπου και θα προκύπτει η τελική έξοδος του φίλτρου με ενεργοποίηση (Σχήμα 4.3.1). Στην είσοδο του κυκλώματος της ReLU (αντί για την έξοδο των πολλαπλασιαστών) εφαρμόζεται και το bias ρεύμα μέσω μιας ελεγχόμενης πηγής ρεύματος καθώς έτσι γίνεται πιο εύκολη η διαδικασία αντιστάθμισης του σταθερού ρεύματος που προσθέτει το κύκλωμα της συνάρτησης ενεργοποίησης στην έξοδο (περισσότερα στην υποενότητα 5.2.2). Η διαφοροποίηση με το αντίστοιχο φίλτρο ενός καναλιού, όπου η πράξη της συνέλιξης και η εφαρμογή της συνάρτησης ενεργοποίησης γίνεται στο ίδιο κύκλωμα, οφείλεται στο γεγονός ότι από προσομοιώσεις διαφαίνεται πως για το



Σχήμα 4.2.2: Παράδειγμα υπολογισμών κυκλώματος 3×3 συνελικτικού φίλτρου ενός καναλιού

συνελικτικό φίλτρο τριών καναλιών, στην συγκεκριμένη εφαρμογή, το υπολογιστικό σφάλμα που εισάγεται αν τα δύο κυκλώματα υλοποιηθούν μαζί είναι σημαντικό. Για τον υπολογισμό της εξόδου του συγκεκριμένου φίλτρου, αυτό «σύρεται» επάνω από την είσοδο τριών καναλιών (Σχήμα 4.3.2) και οι έξοδοι αποθηκεύονται σε αναλογική μνήμη. Επίσης ισχύουν τα ίδια αναφορικά με την δυνατότητα παραλληλοποίησης με το προαναφερθέν φίλτρο ενός καναλιού εισόδου.

4.4 Υλοποίηση πρώτου επιπέδου με ReLU ενεργοποίηση και average pooling

Όπως αναφέρθηκε παραπάνω, μετά το πρώτο συνελικτικό επίπεδο παρεμβάλλεται ένα επίπεδο ομαδοποίησης average pooling. Η υλοποίηση αυτών των επιπέδων θα μπορούσε να γίνει με ξεχωριστά χυχλώματα τα οποία θα λειτουργούσαν ανεξάρτητα. Συγκεκριμένα, θα ήταν δυνατόν να χρησιμοποιηθεί το κύκλωμα της ενότητας 4.2 για να υπολογιστεί το αποτέλεσμα της συνέλιξης με ενεργοποίηση ReLU για το εκάστοτε φίλτρο, αυτό να αποθηκευτεί σε μία αναλογική μνήμη και στη συνέχεια να διαβαστούν αυτά τα αποτελέσματα και να εκτελεστεί τότε η διαδικασία του average pooling. Όμως μπορούμε να αποφύγουμε την επιπλέον κατανάλωση, καθυστέρηση και το σφάλμα που εισάγει η χρήση αναλογικής μνήμης χρησιμοποιώντας τέσσερα χυχλώματα 3×3 φίλτρων δισδιάστατης συνέλιξης (ενότητα 4.2) των οποίων η έξοδος θα οδηγεί σε έναν χοινό χόμβο χαι αχολούθως σε ένα χύχλωμα χλιμάχωσης εξόδου (ενότητα 3.2.2), όπως φαίνεται στο Σχήμα 4.4.1. Η χρήση του χυχλώματος χλιμάχωσης εξόδου γίνεται ώστε τα ρεύματα εισόδου όλων των επιπέδων του διχτύου να μην ξεπερνούν την τιμή των (περίπου) 10nA με σχοπό τον περιορισμό της χατανάλωσης ισχύος από πολύ μεγάλα βάρη αχμών. Αυτό είναι δυνατόν λόγω του γεγονότος ότι η τελική έξοδος του δικτύου προκύπτει από την εφαρμογή του τελεστή argmax στην έξοδο του 5^{ov} στρώματος και έτσι η κλιμάκωση των εξόδων ενός επιπέδου κατά έναν σταθερό συντελεστή δεν θα επηρεάζει το τελικό αποτέλεσμα.

Για τον υπολογισμό τώρα της εξόδου, το κάθε ένα από τα τέσσερα κυκλώματα 3×3 συνελικτικών φίλτρων «σύρονται» πάνω από την εικόνα εισόδου με βήμα 2 ξεκινώντας το καθένα με ένα διαφορετικό offset. Συγκεκριμένα για μία δεδομένη θέση του 1^{ou} (πάνω αριστερά) φίλτρου, το 2^o (πάνω δεξιά) θα βρίσκεται ένα pixel δεξιότερα, το 3^o (κάτω αριστερά) ένα pixel πιο κάτω και το 4^o (κάτω δεξιά) θα βρίσκεται ένα pixel δεξιότερα, το 3^o (κάτω αριστερά) ένα pixel πιο κάτω και το 4^o (κάτω δεξιά) θα βρίσκεται ένα pixel δεξιότερα ποθηκεύεται σε αναλογική μνήμη ενώ ισχύουν τα ίδια περί παραλληλοποίησης με τις παραπάνω υλοποιήσεις.

4.5 Υλοποίηση επιπέδων 2,3,4 με ReLU ενεργοποίηση

Όπως προαναφέρθηκε, τα επίπεδα 2 και 3 περιλαμβάνουν τρία 3×3 συνελικτικά φίλτρα τριών καναλιών εισόδου και το 4^{o} επίπεδο περιλαμβάνει ένα τέτοιο φίλτρο. Έτσι όλα αυτά τα επίπεδα μπορούν να υλοποιηθούν με ένα μόνο κύκλωμα 3×3 συνελικτικού φίλτρου τριών καναλιών (ενότητα 4.3), ένα κύκλωμα ReLU (ενότητα 3.3) και ένα κύκλωμα κλιμάκωσης ρεύματος εξόδου (υποενότητα 3.2.2). Όπως και στην περίπτωση του 1^{ou} επιπέδου, το κύκλωμα κλιμάκωσης ρεύματος εξόδου χρησιμοποιείται ώστε τα ρεύματα εισόδου όλων των επιπέδων του δικτύου να μην ξεπερνούν την τιμή των περίπου 10nA και έτσι να αποφεύγεται η αλόγιστη κατανάλωση ισχύος.

4.6 Υλοποίηση πλήρως συνδεμένου επιπέδου κεφαλής ταξινόμησης

Το πλήρως συνδεδεμένο επίπεδο κεφαλής ταξινόμησης, δεδομένου ότι έχει μέγεθος εισόδου 49 και μέγεθος εξόδου 43, θα περιλαμβάνει 49 × 43 = 2107 ακμές βαρών και 43 ακμές για biases. Η πιο απλή υλοποίηση λοιπόν αυτού του επιπέδου θα ήταν με την χρήση 2107 αναλογικών πολλαπλασιαστών και 43 ηλεκτρονικά ελεγχόμενων πηγών ρεύματος. Κάτι τέτοιο όμως θα οδηγούσε σε μεγάλη κατανάλωση ισχύος και επιφάνεια κυκλώματος. Μια πιο αποδοτική υλοποίηση, η οποία χρησιμοποιείται στην παρούσα εργασία, είναι να υπολογίζεται κάθε μία από τις 43 εξόδους του πλήρως συνδεδεμένου στρώματος σειριακά χρησιμοποιώντας το ίδιο κύκλωμα με 49 αναλογικούς πολλαπλασιαστές και μία ηλεκτρονικά ελεγχόμενη πηγή ρεύματος. Συγκεκριμένα, σε κάθε ένα από τα 43 βήματα που απαιτούνται, οι κανονικοποιημένοι συντελεστές πολλαπλασιασμού των 43 αναλογικών πολλαπλασιαστών και το ρεύμα της πηγή ρεύματος του bias ορίζονται ίσοι με τα βάρη των ακμών που οδηγούν στην εκάστοτε έξοδο. Για το τελικό αυτό επίπεδο του δικτύου η έξοδος ρεύματος δεν κλιμαχώνεται καθώς το ρεύμα αυτό συμβάλει σε αμελητέο βαθμό στην συνολική κατανάλωση ισχύος του εν λόγω συνελικτικού νευρωνικού δικτύου.



Σχήμα 4.3.1: Σχεδιάγραμμα κυκλώματος 3×3 συνελικτικού φίλτρου τριών καναλιών με ReLU ενεργοποίηση



Σχήμα 4.3.2: Παράδειγμα υπολογισμών κυκλώματος 3 × 3 συνελικτικού φίλτρου τριών καναλιού

Σχήμα 4.4.1: Σχεδιάγραμμα κυκλώματος πρώτου επιπέδου με ReLUενεργοποίηση και average pooling





Σχήμα 4.6.1: Σχεδιάγραμμα χυχλώματος πλήρως συνδεδεμένου επιπέδου χεφαλής

Κεφάλαιο 5

Παράδειγμα πραγματικής εφαρμογής

5.1 Σύνολο δεδομένων GTSRB

To GTSRB (German Traffic Sign Recognition Benchmark) [51] ήταν μία δοχιμασία ταξινόμησης πολλαπλών κλάσεων που διεξήχθη στο International Joint Conference on Neural Networks (IJCNN) το 2011. Συγκεκριμένα αυτή αφορούσε την ανάπτυξη μοντέλων μηχανικής μάθησης για την ταξινόμηση πραγματικών εικόνων σημάτων οδικής κυκλοφορίας από το Γερμανικό οδικό δίκτυο σε κάθε μία από 43 πιθανές κλάσεις. Το ομώνυμο σύνολο δεδομένων που χρησιμοποιήθηκε περιλαμβάνει χιλιάδες τέτοιες εικόνες σε διαφορετικές αναλύσεις και συνθήκες φωτισμού. Για την εφαρμογή της εν λόγω διπλωματικής χρησιμοποιώ αυτές τις εικόνες αλλά με σταθερή ανάλυση 28 × 28 και για το πρώτο συνελικτικό επίπεδο του δικτύου θεωρώ ως είσοδο το μέσο όρο των τριών καναλιών τους. Έτσι με αυτούς του περιορισμούς μπορώ να προσομοιώσω ένα πραγματικό σενάριο λειτουργίας του μοντέλου όπου το πρώτο επίπεδο δέχεται ως είσοδο τα ρεύματα εξόδου των pixel μίας μονοχρωματικης κάμερας με ανάλυση 28 × 28.

5.2 Εκπαίδευση μοντέλου

Η εκπαίδευση του μοντέλου του εν λόγω αναλογικού ταξινομητή περιλαμβάνει δυο διακριτές διαδικασίες: Την εκπαίδευση του συνελικτικού νευρωνικού δικτύου σε επίπεδο λογισμικού και την προσαρμογή των ρευμάτων bias και των τάσεων Vsm και Vdm στα κυκλώματα των διαφόρων επιπέδων του ώστε να ελαχιστοποιείται το τελικό σφάλμα.
5.2.1 Εκπαίδευση μοντέλου σε επίπεδο λογισμικού

Η εχπαίδευση του μοντέλου πραγματοποιείται στην γλώσσα προγραμματισμού Ρython χρησιμοποιώντας ένα υποσύνολο 26640 ειχόνων (training dataset) του συνόλου δεδομένων GTSRB. Επειδή αυτό το σύνολο δεδομένων εκπαίδευσης είναι σχετικά μικρό, για να επιτύχω καλύτερη επίδοση του μοντέλου μου χρησιμοποιώ την τεχνική της επαύξησης δεδομένων (data augmentation). Συγκεκριμένα η τεχνική αυτή αναφέρεται στην δημιουργία νέων παραδειγμάτων εχπαίδευσης ενός μοντέλου από τα ήδη υπάρχοντα μέσω ενός μετασχηματισμού ο οποίος δεν αλλάζει τα θεμελιώδη χαραχτηριστικά τους [1]. Έτσι παρέχουμε σε ένα μοντέλο «νέα» δεδομένα τα οποία μπορούν να αυξήσουν την αχρίβεια του και να περιορίσουν το overfitting. Για δεδομένα ειχόνων, συνήθεις μετασχηματισμοί που χρησιμοποιούνται για την επαύξηση δεδομένων είναι η περιστροφή, ο χαθρεπτισμός, η προσθήχη τυχαίου θορύβου χ.α. Στην συγχεχριμένη εφαρμογή χάνω χρήση των μετασχηματισμών της περιστροφής μικρών γωνιών (< 15deg) και της προσθήκης γκαουσιανού θορύβου χαμηλού πλάτους. Προφανώς οι μετασχηματισμένες αυτές ειχόνες εξαχολουθούν να ανήχουν στις εχάστοτε χλάσεις, δεδομένου ότι π.χ. ένα ελαφρώς περιστρεμμένο χαι βρώμιχο σήμα STOP εξαχολουθεί να είναι μια ρεαλιστική ειχόνα ενός σήματος STOP που ένας οδηγός μπορεί να συναντήσει στον πραγματικό κόσμο.

Αχολούθως, χρησιμοποιώντας ως σύνολο δεδομένων εχπαίδευσης το αρχιχό training dataset μαζί με τις μετασχηματισμένες ειχόνες, εχπαιδεύω το συνελιχτιχό νευρωνιχό δίχτυο με την τεχνιχή του mini batching με batch_size = 80. Ως optimizer χρησιμοποιώ τον αλγόριθμο SGDM (stochastic gradient descent with momentum) ο οποίος είναι μια παραλλαγή του αλγορίθμου της χατάβασης δυναμιχού που αναφέρθηχε στο χεφάλαιο 2. Η επιλογή αυτή γίνεται χαθώς ο αλγόριθμος βελτιστοποίησης SGDM συχνά παρέχει μεγαλύτερη ιχανότητα γενίχευσης σε σχέση με πιο σύνθετους optimizers, όπως ο Adam, για εφαρμογές όρασης υπολογιστών [52]. Με αυτήν την διαδιχασία εχπαίδευσης, το παραπάνω συνελιχτιχό νευρωνιχό δίχτυο επιτυγχάνει αχρίβεια ταξινόμησης στο σύνολο δεδομένων ελέγχου (test dataset) (12630 ειχόνων) 88.04%.

5.2.2 Εκπαίδευση μοντέλου σε επίπεδο κυκλωμάτων

Για την παραγωγή των εισόδων του χυχλώματος του πρώτου στρώματος, οι κανονικοποιημένες ($\in [0,1]$)τιμές έντασης των pixel των εικόνων εκπαίδευσης αντιστοιχίζονται σε ρεύματα με εύρος τιμών μεταξύ 0 και 8nA¹. Ακολούθως, οι συντελεστές των εκάστοτε φίλτρων αντιστοιχίζονται στο εύρος τιμών των αναλογικών πολλαπλασιαστών ($\in [-2.1, 2.1]$, με κάποιο περιθώριο ασφαλείας). Στην περίπτωση του πρώτου επιπέδου, η έξοδος του κάθε φίλτρου οδηγείται κατευθείαν στο κύκλωμα υλοποίησης της συνάρτησης ενεργοποίησης ReLU και στη συνέχεια η αθροιστική έξο-

¹Εδώ μπορεί να γίνει η εύλογη υπόθεση ότι αυτό είναι το εύρος τιμών της εξόδου των pixel ενός πραγματικού αισθητήρα.

δος αυτών για τα τέσσερα ίδια φίλτρα οδηγείται σε ένα χύχλωμα χλιμάχωσης εξόδου. Στο σημείο όμως αυτό, αν λάβουμε την έξοδο του συνολιχού αυτού χυχλώματος χαι την χλιμαχώσουμε με τον σωστό συντελεστή παρατηρούνται δύο προβλήματα. Συγχεχριμένα, αν οι τάσεις Vsm χαι Vdm ταυτιστούν με τις τάσεις αρνητιχής χαι θετιχής τροφοδοσίας αντίστοιχα χαι το bias ρεύμα οριστεί ίσο με την θεωρητιχή τιμή, θα παρατηρηθεί μία έξοδος όπως αυτή του Σχήματος 5.2.1, όπου φαίνεται ότι σε αυτήν έχει προστεθεί μια ανεπιθύμητη σταθερή τιμή χαι οι τιμές ρεύματος κοντά στο 0 είναι αλλοιωμένες. Αυτές όμως οι αποχλίσεις μπορούν να διορθωθούν σε μεγάλο βαθμό με τα εξής βήματα: Αρχιχά το δεύτερο πρόβλημα μπορεί να αντιμετωπιστεί ρυθμίζοντας τις τάσεις Vsm χαι Vdm έτσι ώστε οι χαμηλές τιμές ρεύματος να μην αλλοιώνονται και παράλληλα να υπάρχει αρχετό εύρος τάσης ώστε τα μεγάλα ρεύματα να μην συμπιέζονται. Αχολούθως, η σταθερή απόχλιση του ρεύματος εξόδου που παρατηρείται μετά από το παραπάνω βήμα μπορεί να αντισταθμιστεί μεταβάλλοντας κατάλληλα το ρεύμα bias (Σχήμα 5.2.2).

Για τα επόμενα τρία συνελικτικά επίπεδα, δεδομένου ότι η εφαρμογή της συνάρτησης ενεργοποίησης γίνεται ξεχωριστά από τον υπολογισμό του αποτελέσματος κάθε φίλτρου, δεν παρατηρείται αλλοίωση των χαμηλών τιμών ρεύματος. Έτσι, οι τιμές των τάσεων Vsm και Vdm για το ξεχωριστό κύκλωμα ReLU ενεργοποίησης ορίζονται ίσες με αυτές των Vss και Vdd αντίστοιχα. Παρατηρείται όμως θετικό DC offset το οποίο αντισταθμίζεται πάλι μεταβάλλοντας το ρεύμα bias, το οποίο εφαρμόζεται στην είσοδο του κυκλώματος ενεργοποίησης αντί για την έξοδο του εκάστοτε φίλτρου.

Αντίστοιχα με τα παραπάνω, το θετικό DC offset που παρατηρείται και στο τελευταίο πλήρως συνδεδεμένο στρώμα του εν λόγω CNN αντισταθμίζεται πάλι με την κατάλληλη προσαρμογή των ρευμάτων bias. Σε αυτό όμως το επίπεδο πρέπει παράλληλα να εξασφαλιστεί ότι για σχεδόν όλες τις πιθανές εισόδους του μοντέλου, η έξοδος του πλήρως συνδεδεμένου επιπέδου θα περιλαμβάνει τουλάχιστον μία θετική τιμή ρεύματος. Αυτό είναι απαραίτητο επειδή το κύκλωμα Winner-Take-All, που υλοποιεί τον τελεστή argmax, είναι σχεδιασμένο για να συγκρίνει ορθά μόνο θετικές τιμές ρεύματος και στην περίπτωση αποκλειστικά αρνητικών εισόδων αυτό δεν θα παρήγαγε την επιθυμητή έξοδο. Έτσι σε περίπτωση που αυτή η απαίτηση δεν εξασφαλίζεται, προστίθεται στην έξοδο μια ακόμη σταθερή τιμή ρεύματος η οποία προφανώς δεν θα επηρεάζει το αποτέλεσμα του τελεστή argmax. Για την εν λόγω όμως εφαρμογή, αυτή η απαίτηση εξασφαλίζεται χωρίς την προσθήκη επιπλέον ρεύματος.

Αχόμη, για όλα τα συνελιχτικά επίπεδα, μετά από τις προαναφερθέντες προσαρμογές, ορίζεται η τιμή της τάσης εισόδου του χυχλώματος χλιμάχωσης εξόδου ώστε κάθε μία από τις τελικές εξόδους τους να έχει μέγιστη τιμή ρεύματος στο εύρος [8nA, 10nA]. Το βήμα αυτό σε συνδυασμό με την γραμμική αντιστοίχηση των τιμών εισόδων και των βαρών των αχμών του διχτύου με τις αντίστοιχες χυχλωματικές τιμές ισοδυναμεί αριθμητικά με την χλιμάχωση της εξόδου του πλήρως συνδεμένου επιπέδου του πρωτότυπου μοντέλου με μία σταθερά. Αυτή όμως η κλιμάχωση δεν έχει επίπτωση στην τελική έξοδο ταξινόμησης δεδομένου ότι αυτή παράγεται μέσω του τελεστή argmax.

Όλες οι παραπάνω διαδικασίες γίνονται τροφοδοτώντας στο μοντέλο 100 εικόνες από το σύνολο δεδομένων εκπαίδευσης και προσαρμόζοντας τις εκάστοτε κυκλωματικές παραμέτρους κατάλληλα ώστε να παρατηρείται η μέγιστη δυνατή ταύτιση μεταξύ της κυκλωματικής εξόδου του εκάστοτε επιπέδου, μετά από την ορθή κλιμάκωση, και της εξόδου του πρωτότυπου μοντέλου. Γενικά οι σταθερές αποκλίσεις που παρατηρούνται στην έξοδο των κυκλωμάτων, χωρίς αντιστάθμιση, για όλα τα στρώματα, που αποτελεί το σημαντικότερο εκ των δύο προαναφερθέντων προβλημάτων, οφείλεται στην μικρή διαφορά της συνάρτησης μεταφοράς των αναλογικών πολλαπλασιαστών για θετικούς και αρνητικούς συντελεστές και, κυρίως, στους καθρέπτες ρεύματος της ReLU. Αξίζει όμως να σημειωθεί πως η ίδια διαδικασία της τροφοδότησης κάποιων δοκιμαστικών εισόδων στο δίκτυο, της εξαγωγής των εξόδων των διαφόρων επιπέδων και της ανάλογης προσαρμογής των ηλεκτρονικά ελεγχόμενων κυκλωματικών παραμέτρων μπορεί να χρησιμοποιηθεί για την αντιστάθμιση τέτοιων σφαλμάτων ανεξαρτήτως αιτίας. Έτσι, με αυτόν τον τρόπο μπορούν να αντιμετωπιστούν, σε κάποιο βαθμό, και σφάλματα PVT.

5.3 Αποτελέσματα ταξινόμησης συνόλου δεδομένων GTSRB

Για την αξιολόγηση της αχρίβειας του παραπάνω χυχλωματιχού μοντέλου, η έξοδος του συγχρίνεται με αυτήν του πρωτότυπου μοντέλου σε επίπεδο λογισμιχού για 1000 τυχαία επιλεγμένα δείγματα από το σύνολο δεδομένων αξιολόγησης. Συγχεχριμένα, για τις χυχλωματιχές παραμέτρους που ορίστηχαν χατά την διαδιχασία εχπαίδευσης, τροφοδοτώ στα χυχλώματα των διαφόρων επιπέδων τέσσερα διαφορετιχά batches 250 δειγμάτων από τα οποία προχύπτουν τα αποτελέσματα του Πίναχα 5.1.

Εδώ να σημειωθεί πως για την αχρίβεια του χυχλωματιχού μοντέλου λαμβάνονται υπόψη μόνο τα δείγματα που ταξινομεί σωστά αυτό αλλά και το πρωτότυπο μοντέλο. Επίσης δεν προσμετρώνται τα ορθά ταξινομημένα δείγματα για τα οποία η έξοδος της νικήτριας κλάσης στο κύκλωμα Winner-Take-All είναι κάτω από 5nA καθώς σε αυτήν την περίπτωση, σε πραγματικές συνθήκες η νικήτρια κλάση θα ήταν πρακτικά τυχαία.

Από τα αποτελέσματα της προσομοίωσης του εν λογω αναλογικού συνελτικτικού δικτύου παρατηρώ πως το σφάλμα ταξινόμησης μεταξύ αυτού και του πρωτότυπου δικτύου είναι ιδιαίτερα μικρό ενώ. Για την παρούσα προσομοίωση όλα τα κυκλώματα των διαφόρων επιπέδων λειτούργησαν σε συχνότητα 100kHz και έτσι, αγνοώντας τις καθυστερήσεις λόγω μνήμης, το αναλογικό μοντέλο αυτό παρήγαγε την έξοδο για κάθε εικόνα σε $(3 \cdot 13 \cdot 13 + 6 \cdot 11 \cdot 11 + 6 \cdot 9 \cdot 9 + 2 \cdot 7 \cdot 7 + 43) \cdot 10 \mu s = 18.6ms$. Παράλληλα η μέγιστη κατανάλωση των διαφόρων επιπέδων όπως μετρήθηκε για 250

	Software model	Circuit model
Batch 1	91.2%	89.2%
Batch 2	86.0%	82.4%
Batch 3	88.8%	88.0%
Batch 4	87.2%	85.2%
Average	88.3%	86.2%

Πίναχας 5.1: Αχρίβεια ταξινόμησης 1000 δειγμάτων του συνόλου δεδομένων GTSRB

δείγματα είναι εξαιρετικά μικρή (Πίνακας 5.2) επιβεβαιώνοντας έτσι την μεγάλη πρακτική χρησιμότητα της προτεινόμενης υλοποίησης.

Πίνα
χας 5.2: Μέγιστη κατανάλωση ενέργειας μετρημένη για 250 δείγματα του συνό
λου δεδομένων GTSRB

	Peak Power (μW)
First layer circuit	1.77
Second layer circuit (without ReLU)	0.743
ReLU circuit (for second layer)	0.061
Third layer circuit (without ReLU)	0.742
ReLU circuit (for third layer)	0.063
Fourth layer circuit (without ReLU)	0.787
ReLU circuit (for fourth layer)	0.062
Fully connected layer circuit	1.38
WTA circuit	0.144

Σχήμα 5.2.1: Έξοδος ρεύματος επιπέδου average pooling για το δεύτερο και το τρίτο φίλτρο του πρώτου συνελικτικού επιπέδου χωρίς προσαρμογή των Ibias, Vsm, Vdm





Σχήμα 5.2.2: Έξοδος ρεύματος επιπέδου average pooling για το δεύτερο και το τρίτο φίλτρο του πρώτου συνελικτικού επιπέδου μετά απο προσαρμογή των Ibias, Vsm, Vdm





Κεφάλαιο 6

Συμπεράσματα και μελλοντική δουλειά

Στην παρούσα διπλωματική εργασία παρουσιάστηκε μια πλήρως αναλογική υλοποίηση ση εξαιρετικά χαμηλής κατανάλωσης ενός βαθέως συνελικτικού νευρωνικού δικτύου. Συγκεκριμένα αναλύθηκε η σχεδίαση κυκλωμάτων για την υλοποίηση συνελικτικών επιπέδων διαφορετικών διαστάσεων και του πλήρως συνδεδεμένου στρώματος του δικτύου καθώς και για την υλοποίηση της ReLU συνάρτηση ενεργοποίησης και του τελεστή argmax. Το αναλογικό αυτό μοντέλο κατάφερε σε ένα δείγμα 1000 εισόδων να επιτύχει ακρίβεια 86.2% η οποία απέχει μόλις 2.1% από αυτήν του πρωτότυπου μοντέλου σε επίπεδο λογισμικού με την μέγιστη κατανάλωση ισχύος του να κυμαίνεται κάτω από $1.8\mu W$. Η υλοποίηση των επιπέδων με εκπαιδεύσιμες παραμέτρους έγινε με την χρήση ενός πρωτότυπου πλήρως ηλεκτρονικά ελεγχόμενου κυκλώματος αναλογικού πολλαπλασιασμού το οποίο μπορεί να λειτουργήσει με υψηλή ακρίβεια για εξαιρετικά χαμηλά ρεύματα εισόδου (< 1nA) και για ένα ευρύ φάσμα εφαρμογών.

Στην παρούσα εργασία επιλέχθηκε η υλοποίηση ενός σχετικά απλού μοντέλου για την ταξινόμηση εικόνων οδικών σημάτων ώστε να φανεί ξεκάθαρα το πλεονέκτημα της εξαιρετικά χαμηλής κατανάλωσης και υψηλής ακρίβειας που μπορεί να προσφέρει η προτεινόμενη αρχιτεκτονική για ένα πραγματικό πρόβλημα. Το γεγονός όμως ότι τα κυκλώματα που χρησιμοποιήθηκαν για την υλοποίηση των επιπέδων αυτού του μοντέλου είναι πλήρως ηλεκτρονικά ελεγχόμενα δίνει την δυνατότητα αξιοποίησης τους για την υλοποίηση διαφορετικών νευρωνικών δικτύων που χρησιμοποιούν για τα επίπεδα τους τα ίδια δομικά στοιχεία με το προτεινόμενο μοντέλο. Ακόμη, με βάση τα πρωτότυπα κυκλώματα αναλογικών πολλαπλασιαστών που αυτή η εργασία εισάγει, είναι δυνατόν να κατασκευαστούν, με τις ίδιες σχεδιαστικές αρχές, και πιο σύνθετα μοντέλα νευρωνικών δικτύων. Επίσης, με την χρήση διαφορετικών σημάτων εισόδου τάσης για κάθε ένα από τα δύο κυκλώματα σιγμοειδούς των εν λόγω αναλογικών πολλαπλασιαστών είναι δυνατόν να αντισταθμιστούν σε μεγάλο βαθμό σφάλματα λόγω transistor mismatches ή άλλων παραγόντων. Έτσι σε συνδυασμό με την αυτοματοποίηση της διαδικασίας διόρθωσης των ανεπιθύμητων σφαλμάτων λόγω DC offset μπορούν να σχεδιαστούν αναλογικοί επεξεργαστές οι οποίοι με έναν σχετικά μικρό αριθμό υπολογιστικών μονάδων εξαιρετικά χαμηλής κατανάλωσης και κάποια κελιά αναλογικής μνήμης θα μπορούν να υλοποιούν πολλά διαφορετικά μοντέλα σύνθετων νευρωνικών δικτύων προσομοιάζοντας την λειτουργία ενός αναλογικού FPGA.

Bibliography

- [1] C. C. Aggarwal, Neural Networks and Deep Learning. Springer, 2018.
- [2] G. Rätsch, "A brief introduction into machine learning," Friedrich Miescher Laboratory of the Max Planck Society, pp. 1–6, 2004.
- [3] S. Bhattacharyya and N. Dey, Trends in Deep Learning Methodologies. ACA-DEMIC PRESS, 2021.
- [4] I. H. Sarker, "Machine learning: Algorithms, real-world applications and research directions," SN computer science, vol. 2, no. 3, p. 160, 2021.
- [5] V. Sze, Y.-H. Chen, J. Emer, A. Suleiman, and Z. Zhang, "Hardware for machine learning: Challenges and opportunities," in 2017 IEEE custom integrated circuits conference (CICC), pp. 1–8, IEEE, 2017.
- [6] I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. MIT Press, 2016.
- [7] M. I. Jordan and T. M. Mitchell, "Machine learning: Trends, perspectives, and prospects," *Science*, vol. 349, no. 6245, pp. 255–260, 2015.
- [8] R. S. Sutton and A. G. Barto, *Reinforcement Learning: An Introduction*. MIT Press, 2015.
- [9] N. $K\omega\nu\sigma\tau\alpha\nu\tau\iota\nu\alpha$, $\Pi\rho\sigma\sigma\sigma\mu\iota\omega\sigma\eta \Phi\nu\sigma\iota\sigma\lambda\sigma\gamma\iota\kappa\omega\nu \Sigma\nu\sigma\tau\eta\mu\dot{\alpha}\tau\omega\nu$. Tziolas Publications, 2016.
- [10] A. Prieto, B. Prieto, E. M. Ortigosa, E. Ros, F. Pelayo, J. Ortega, and I. Rojas, "Neural networks: An overview of early research, current frameworks and new challenges," *Neurocomputing*, vol. 214, pp. 242–268, 2016.
- [11] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Imagenet classification with deep convolutional neural networks," *Advances in neural information pro*cessing systems, vol. 25, 2012.
- [12] C. M. Bishop, Pattern recognition and machine learning. Springer, 2006.

- [13] D. H. Hubel, T. N. Wiesel, et al., "Receptive fields of single neurones in the cat's striate cortex," J physiol, vol. 148, no. 3, pp. 574–591, 1959.
- [14] Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, "Gradient-based learning applied to document recognition," *Proceedings of the IEEE*, vol. 86, no. 11, pp. 2278–2324, 1998.
- [15] K. Simonyan and A. Zisserman, "Very deep convolutional networks for largescale image recognition," arXiv preprint arXiv:1409.1556, 2014.
- [16] V. Alimisis, G. Gennis, K. Touloupas, C. Dimas, N. Uzunoglu, and P. P. Sotiriadis, "Nanopower integrated gaussian mixture model classifier for epileptic seizure prediction," *Bioengineering*, vol. 9, no. 4, p. 160, 2022.
- [17] V. Alimisis, G. Gennis, C. Dimas, and P. P. Sotiriadis, "An analog bayesian classifier implementation, for thyroid disease detection, based on a lowpower, current-mode gaussian function circuit," in 2021 International conference on microelectronics (ICM), pp. 153–156, IEEE, 2021.
- [18] K. H. Lee and N. Verma, "A low-power processor with configurable embedded machine-learning accelerators for high-order and adaptive analysis of medical-sensor signals," *IEEE Journal of Solid-State Circuits*, vol. 48, no. 7, pp. 1625–1637, 2013.
- [19] E. Nurvitadhi, J. Sim, D. Sheffield, A. Mishra, S. Krishnan, and D. Marr, "Accelerating recurrent neural networks in analytics servers: Comparison of fpga, cpu, gpu, and asic," in 2016 26th International Conference on Field Programmable Logic and Applications (FPL), pp. 1–4, IEEE, 2016.
- [20] C. Farabet, Y. LeCun, K. Kavukcuoglu, E. Culurciello, B. Martini, P. Akselrod, and S. Talay, "Large-scale fpga-based convolutional networks," *Scaling up machine learning: parallel and distributed approaches*, vol. 13, no. 3, pp. 399–419, 2011.
- [21] P. Marwedel, Embedded System Design. Springer, 2006.
- [22] C. Demirkiran, L. Nair, D. Bunandar, and A. Joshi, "A blueprint for precise and fault-tolerant analog neural networks," *Nature Communications*, vol. 15, no. 1, p. 5098, 2024.
- [23] C. Silvano, D. Ielmini, F. Ferrandi, L. Fiorin, S. Curzel, L. Benini, F. Conti, A. Garofalo, C. Zambelli, E. Calore, *et al.*, "A survey on deep learning hardware accelerators for heterogeneous hpc platforms," *arXiv preprint arXiv:2306.15552*, 2023.

- [24] W. Haensch, T. Gokmen, and R. Puri, "The next generation of deep learning hardware: Analog computing," *Proceedings of the IEEE*, vol. 107, no. 1, pp. 108–122, 2018.
- [25] J.-s. Seo, J. Saikia, J. Meng, W. He, H.-s. Suh, Y. Liao, A. Hasssan, I. Yeo, et al., "Digital versus analog artificial intelligence accelerators: Advances, trends, and emerging designs," *IEEE Solid-State Circuits Magazine*, vol. 14, no. 3, pp. 65–79, 2022.
- [26] V. Alimisis, M. Gourdouparis, G. Gennis, C. Dimas, and P. P. Sotiriadis, "Analog gaussian function circuit: Architectures, operating principles and applications," *Electronics*, vol. 10, no. 20, p. 2530, 2021.
- [27] I. Dadras, M. H. Ahmadilivani, S. Banerji, J. Raik, and A. Abloo, "An efficient analog convolutional neural network hardware accelerator enabled by a novel memoryless architecture for insect-sized robots," in 2022 11th International Conference on Modern Circuits and Systems Technologies (MO-CAST), pp. 1–6, IEEE, 2022.
- [28] M. Chao, All Analog CNN Accelerator with RRAMs for Fast Inference. PhD thesis, Massachusetts Institute of Technology, 2022.
- [29] V. Alimisis, G. Gennis, M. Gourdouparis, C. Dimas, and P. P. Sotiriadis, "A low-power analog integrated implementation of the support vector machine algorithm with on-chip learning tested on a bearing fault application," *Sensors*, vol. 23, no. 8, p. 3978, 2023.
- [30] S. Sadasivuni, S. P. Bhanushali, I. Banerjee, and A. Sanyal, "In-sensor neural network for high energy efficiency analog-to-information conversion," *Scientific reports*, vol. 12, no. 1, p. 18253, 2022.
- [31] H. Valavi, P. J. Ramadge, E. Nestler, and N. Verma, "A mixed-signal binarized convolutional-neural-network accelerator integrating dense weight storage and multiplication for reduced data movement," in 2018 IEEE Symposium on VLSI Circuits, pp. 141–142, IEEE, 2018.
- [32] Z. Jiang, S. Yin, J.-S. Seo, and M. Seok, "C3sram: An in-memory-computing sram macro based on robust capacitive coupling computing mechanism," *IEEE Journal of Solid-State Circuits*, vol. 55, no. 7, pp. 1888–1897, 2020.
- [33] Z. Chen, X. Chen, and J. Gu, "15.3 a 65nm 3t dynamic analog ram-based computing-in-memory macro and cnn accelerator with retention enhancement, adaptive analog sparsity and 44tops/w system energy efficiency," in 2021 IEEE International Solid-State Circuits Conference (ISSCC), vol. 64, pp. 240–242, IEEE, 2021.

- [34] M. N. Bojnordi and E. Ipek, "Memristive boltzmann machine: A hardware accelerator for combinatorial optimization and deep learning," in 2016 IEEE International Symposium on High Performance Computer Architecture (HPCA), pp. 1–13, IEEE, 2016.
- [35] P. Chi, S. Li, C. Xu, T. Zhang, J. Zhao, Y. Liu, Y. Wang, and Y. Xie, "Prime: A novel processing-in-memory architecture for neural network computation in reram-based main memory," ACM SIGARCH Computer Architecture News, vol. 44, no. 3, pp. 27–39, 2016.
- [36] M. S. Asghar, S. Arslan, and H. Kim, "Analog convolutional operator circuit for low-power mixed-signal cnn processing chip," *Sensors*, vol. 23, no. 23, p. 9612, 2023.
- [37] F. B. Gencer, X. Xhafa, B. B. İnam, and M. B. Yelten, "Design and validation of an artificial neural network based on analog circuits," *Analog Integrated Circuits and Signal Processing*, vol. 106, pp. 475–483, 2021.
- [38] V. Alimisis, A. Papathanasiou, E. Georgakilas, N. P. Eleftheriou, and P. P. Sotiriadis, "An ultra-low power adjustable current-mode analog integrated general purpose artificial neural network classifier," *AEU-International Journal of Electronics and Communications*, vol. 186, p. 155467, 2024.
- [39] F. Rosenblatt, "The perceptron: a probabilistic model for information storage and organization in the brain.," *Psychological review*, vol. 65, no. 6, p. 386, 1958.
- [40] M. Grove and J. Blinkhorn, "Neural networks differentiate between middle and later stone age lithic assemblages in eastern africa," *PloS one*, vol. 15, no. 8, p. e0237528, 2020.
- [41] S. Shalev-Shwartz and S. Ben-David, Understanding machine learning: From theory to algorithms. Cambridge university press, 2014.
- [42] S. Haykin, Neural networks and learning machines, 3/E. Pearson Education, 2009.
- [43] H. T. Nguyen, N. R. Prasad, C. L. Walker, and E. A. Walker, A first course in fuzzy and neural control. Chapman and Hall/CRC, 2002.
- [44] V. Nair and G. E. Hinton, "Rectified linear units improve restricted boltzmann machines," in *Proceedings of the 27th international conference on machine learning (ICML-10)*, pp. 807–814, 2010.
- [45] A. Sedra and K. Smith, *Microelectronic circuits 8th edition*. Oxford University Press, 2015.

- [46] B. Razavi, Design of analog CMOS integrated circuits. McGraw-Hill Education, 2017.
- [47] P. R. Gray, P. J. Hurst, S. H. Lewis, and R. G. Meyer, Analysis and design of analog integrated circuits. John Wiley & Sons, 2009.
- [48] M. Gourdouparis, V. Alimisis, C. Dimas, and P. P. Sotiriadis, "An ultra-low power,±0.3 v supply, fully-tunable gaussian function circuit architecture for radial-basis functions analog hardware implementation," *AEU-International Journal of Electronics and Communications*, vol. 136, p. 153755, 2021.
- [49] J. Lazzaro, S. Ryckebusch, M. A. Mahowald, and C. A. Mead, "Winner-takeall networks of o (n) complexity," Advances in neural information processing systems, vol. 1, 1988.
- [50] L. Zhao and Z. Zhang, "A improved pooling method for convolutional neural networks," *Scientific Reports*, vol. 14, no. 1, p. 1589, 2024.
- [51] J. Stallkamp, M. Schlipsing, J. Salmen, and C. Igel, "Man vs. computer: Benchmarking machine learning algorithms for traffic sign recognition," *Neural networks*, vol. 32, pp. 323–332, 2012.
- [52] I. Loshchilov, "Decoupled weight decay regularization," arXiv preprint arXiv:1711.05101, 2017.