RESEARCH ARTICLE                                                                OPEN ACCESS

# Feed Forward Pipelined Accumulate Unit for the Machine Learning

K. Darani*, Dr. P. Sukumar**

*(Department of Electronics and Communication Engineering, Nandha Engineering College,Erode,India*
Email: daranikettimuthu20@gmail.com)
** (*Department of Electronics and Communication Engineering, Nandha Engineering College,Erode,India*
Email: Sukumarwin@gmail.com)

----------------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*---------------------------------

## Abstract:

MAC computations represent an immense bit of AI stimulating specialist exercises. The pipelined structure is by and large grasped to improve the execution by diminishing the length of fundamental ways. An extension in the amount of flip-tumbles due to pipelining, in any case, generally achieves an enormous area and force increase. A tremendous number of flip-flops are normally required to meet the feedforward-cutset rule. In perspective on the discernment that this standard can be free in AI applications, we propose a pipelining technique that discards a bit of the flip-flops explicitly. The re-enactment results show that the proposed MAC unit achieved a 20% essentialness saving and a 20% zone decline differentiated and the normal pipelined MAC.

*Keywords*: **Terms—Hardware accelerator, machine learning, multiply–accumulate (MAC) unit, Pipelining.**

----------------------------------------\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*\*---------------------------------

## I. INTRODUCTION

Starting late, the significant neural framework (DNN) rose as a basic resource for various applications including picture course of action and talk affirmation. Since a colossal proportion of vector-network enlargement figuring's are required in a common DNN application, a variety of submitted hardware for AI have been proposed to enliven the computations. In an AI reviving operator, incalculable copy gather (MAC) units are joined for equal counts, and timing-essential methods for the system are routinely found in the unit.

A multiplier normally contains a couple of computational parts including a midway thing age, a fragment development, and a last extension. An authority involves the pass on causing snake. Long essential courses through these stages lead to the introduction debasement of the general structure. To restrain this issue, various methods have been mulled over. The Wallace and Dadda multipliers are remarkable models for the achievement of a snappy area development, and the pass on lookahead (CLA) snake is routinely used to lessen the fundamental path in the aggregator or the last extension period of the multiplier.

It is outstanding that pipelining is one of the most prominent methodologies for expanding the activity clock recurrence. In spite of the fact that pipelining is an effective method to diminish the basic way

delays, it brings about an expansion in the region and the power utilization because of the addition of many flip-flops. Specifically, the quantity of flip-flops will in general be huge on the grounds that the flip-flops must be embedded in the feed forward-cutset to guarantee useful equity when the pipelining. The issue exacerbates as the quantity of pipeline stages is expanded.

The basic idea of this paper is the ability to relax up the feed forward-cutset rule in the MAC structure for AI applications, considering the way that solitary the last worth is used out of the tremendous number of increment assortments. So to speak, not exactly equivalent to the utilization of the customary MAC unit, moderate hoarding characteristics are not used here, and subsequently, they don't ought to be directly as long as the last worth is correct. Under such a condition, the last worth can end up right if each twofold commitment of the adders inside the MAC checks out the figuring once and just once, free of the cycle. Likewise, it isn't critical to characterize an exact pipeline limit.

In perspective on the as of late explained suspected, this paper proposes a feedforward sans cutset (FCF) pipelined MAC plan that is specific for a prevalent AI stimulating operator. The proposed arrangement system reduces the zone and the force usage by lessening the amount of implanted flip-flops for the pipelining.

## II.PRELIMINARY: FEEDFORWARD-CUTSET RULE FOR PIPELINING

It is remarkable that pipelining is perhaps the best ways to deal with reduce the essential route delay, right now the clock repeat. This lessening is cultivated through the incorporation of flip-flops into the datapath. Fig. 1(a) exhibits the square chart

of a three-tap finiteimpulse response (FIR) channel [11] $y[n] = ax[n] + bx[n - 1] + cx[n - 2]$. (1) Fig. 1(b) and (c) shows pipelining models as for the FIR channel. Despite diminishing fundamental path delays through pipelining, it is similarly basic to satisfy useful consistency while pipelining. When the flip-flops are implanted to ensure utilitarian reasonableness is known as the feedforward-cutset. The implications of cutset and feedforward-cutset are according to the accompanying [11]: Cutset: A ton of the edges of a graph with the ultimate objective that if these edges are removed from the outline, and the chart pushes toward getting disengaged. Feedforward-cutset: A cutset where the data move in the forward heading on most of the cutset edges. Fig. 1(b) shows an instance of significant pipelining. The two-sort out pipelined FIR channel is created by embeddings two flip-tumbles along feedforward-cutset. Strikingly, Fig. 1(c) seems an instance of invalid pipelining. Valuable correspondence isn't guaranteed for this circumstance considering the way that the flip-flops are not inserted viably along the feedforward-cutset.
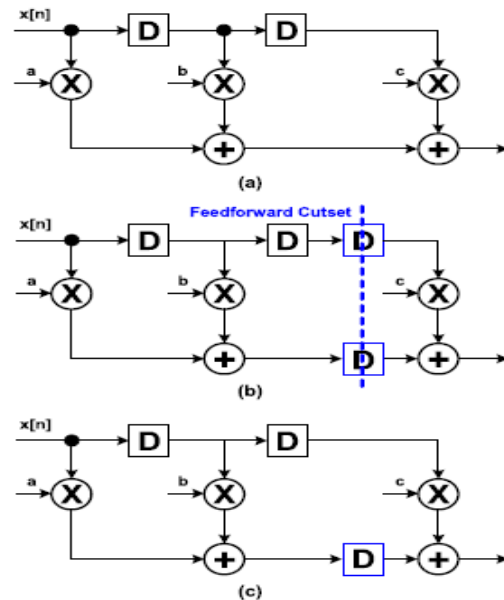

(a)
(b)
(c)

---

FIG.1 Block diagrams of the three-tap FIR filter [11]. (b) Valid pipelining. (c) Invalid pipelining. "D" indicates a flip-flop.

## III. PROPOSED FCF PIPELINING

Fig. 2 shows examples of the two-organize 32-piece pipelined authority (PA) that relies upon the swell pass on snake (RCA). A[31 : 0] addresses data that move from the outside to the information pad register. AReg [31 : 0] addresses the data that are taken care of in the information support. S[31 : 0] addresses the data that are taken care of in the yield bolster register due to the social occasion. In the conventional PA structure [Fig. 2(a)], the flip-flops must be implanted along the feedforward-cutset to ensure valuable equalization. Since the gatherer in Fig. 2(a) contains two pipeline masterminds, the amount of additional flip-flops for the pipelining is 33 (dim toned flip-flops). In case the gatherer is pipelined to the n-orchestrate, the amount of inserted flip-flops winds up 33(n−1), which avows that the amount of flip-flops for the pipelining augmentations through and through as the amount of pipeline stages is extended. Fig. 2(b) exhibits the proposed FCF-PA. For the FCF-PA, only one flip-flop is inserted for the two-stag pipelining. In like manner, the amount of additional flip-flops for the n-sort out pipeline is n − 1 specifically. In the normal PA, the correct assortment estimations of the extensive number of commitments up to the contrasting check cycle are made in each check cycle as showed up in the arranging diagram of Fig. 2(a). A two-cycle qualification exists between the data and the looking at yield on account of the two-mastermind pipeline. On the other hand, in the proposed structure, only the last accumulation result is genuine as showed up in the arranging diagram of Fig. 2(b).

### A. Modified FCF-PA for Further Power Reductions

Notwithstanding the way that the proposed FCF-PA can diminish the zone and the force use by overriding the CLA, there are certain data conditions in which the undesired data change in the yield support occurs, right now the force profitability when 2's enhancement numbers are used. Fig. 3 shows an instance of the undesired data progress. The data sources are 4-piece 2's enhancement combined numbers. AReg[7 : 4] is the sign development of AReg[3], which is the sign bit of AReg[3 : 0]. In the customary pipelining [Fig. 3 (left)], the accumulation result (S) in cycle 3 and the data set aside in the data pad (AReg) in cycle 2 are incorporated and set aside in the yield support (S) in cycle 4. For this circumstance, the "1" in AReg[2] in cycle 2 and the "1" in S[2] in cycle 3 are incorporated, subsequently making a pass on. The pass on is transmitted to the upper part of the S, and along these lines, S[7:4]
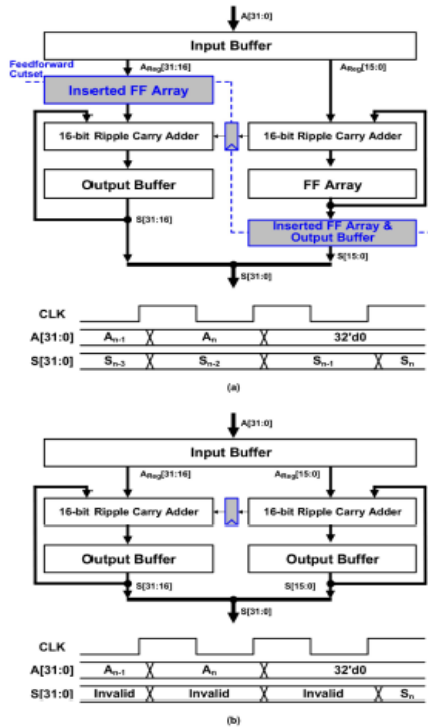
remains as "0000" in cycle 4.



Fig. 2. Schematics and timing diagrams of two-stage 32-bit accumulators.(a) Conventional PA. (b) Proposed FCF-PA.



Fig. 3. Example of an undesired data transition in the two-stage 8-bit Pas with 4-bit 2's complement input numbers.

numbers in each extension organize using the half-adders or conceivably full adders and after that for the demise of the results to the accompanying development orchestrate. Since MAC estimations rely upon such augmentations, the proposed pipelining methodology can in like manner be applied to the AI unequivocal MAC structure. Right now, proposed pipelining technique is applied to the MAC building by using the unique typical for Dadda multiplier. The Dadda multiplier plays out the segment extension thusly to the Wallace multiplier which is extensively used, and it has less zone and shorter fundamental path delay than the Wallace multiplier Fig. 4 shows the pipelined area extension structures in the Dadda multiplier. The Dadda multiplier plays out the segment development to reduce the height of each stage. In case a particular segment starting at now satisfies the target stature for the accompanying section extension sort out, by then no action is performed during the stage Using this property, the proposed pipelining method can be applied to the MAC structure moreover. Fig. 4 (a) will be an instance of pipelining where the customary method is used. Most of the edges in the feedforward-cutset are at risk to pipelining. Of course, in the proposed FCF pipelining case [Fig. 4 (b)], if a center point in the segment development arrange doesn't need to partake in the stature decline, it will in general be banned from the pipelining [the bundle in the spotted box of Fig. 4(b)]. Toward the day's end, in the customary pipelining system, all of the edges in the feedforward-cutset must be pipelined to ensure functional consistency paying little psyche to an arranging slack of each edge [Fig. 4(a)].

## IV. WORKING

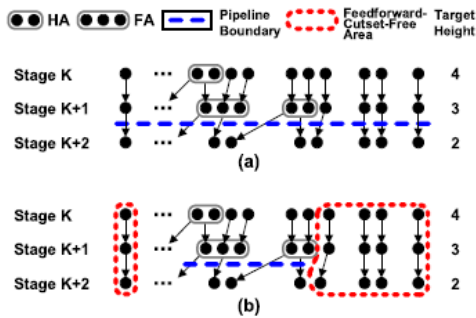The section development in the MAC action is for the figuring of twofold

Fig. 4. Pipelined column addition structure with the Dadda multiplier. (a) Conventional pipelining. (b) Proposed FCF pipelining. HA: half-adder. FA: full adder.
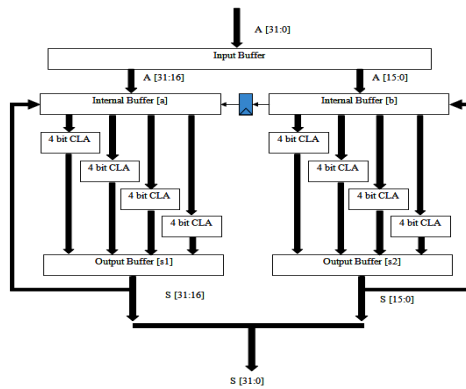


Fig. 5. Block Diagram of partially Pipelined MFCF.

## V. RESULT

We survey the proposed FCF pipelining methodology right now. To begin with, for utilization of the gatherer simply case, twofold weight-frameworks are considered. , we choose the amount of bits in the gatherer to be [16 (InputFeature) +11 (Accumulation) =] 27-piece. For the MAC case, the 16-piece 2's enhancement number is used for both the data feature and weight. Taking everything into account, the amount of bits in the last snake is set out to be [16 × 2 (Multiplication) +11 (Accumulation) =] 43-piece. The arrangement is joined with the passage level cells in a 65-nm CMOS development using Synopsys Design Compiler. For some uncommonly creates,

"set_dont_touch" heading of Design Compiler is used after the prompt dispatch of the cells in the standard library, rather than the association of cells using the register-move level delineation.
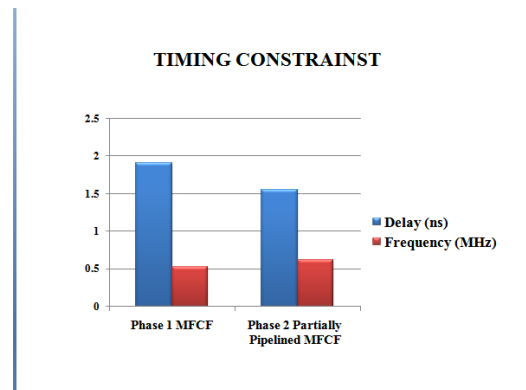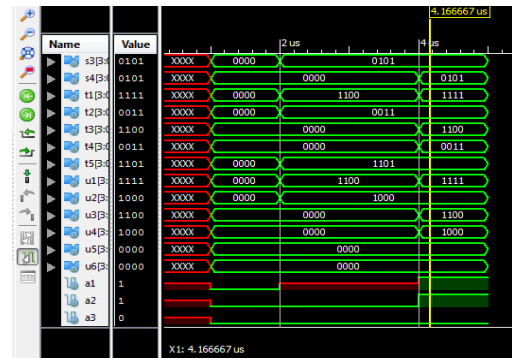


Fig. 6. Comparison of Time



Fig. 7. Simulation Result

For the appraisal of the force use, we run the time touchy examination with a value charge dump record in the PrimeTime PX. Both the genuine data features (ImageNet instructive assortment) and subjective vectors created by the pseudorandom number generator (PRNG) are sustained as information data. We structure the CLA by simply depicting it as "A + B" in Verilog Hardware Description coordinating/streamlining using Design Compiler. All diversion results for the zone and the control use join data and yield pads. The numbers above bars in the figures show institutionalized locale/ability

to the CLA-based aggregator (gatherer simply case) or "Mac + CLA" plan (MAC case).

## VI. CONCLUSION

We exhibited the FCF pipelining technique right now. In the proposed arrangement, the amount of flip-droops in a pipeline can be diminished by loosening up the feedforward-cutset restriction, on account of the exceptional typical for the AI figuring. We applied the FCF pipelining strategy to the aggregator (FCF-PA) structure, and a while later improved the force dispersal of FCF-PA by reducing the chance of undesired data propels (MFCF-PA). The proposed arrangement was in addition broadened, and applied to the MAC unit (FCF-MAC). For the evaluation, the standard and proposed MAC structures were joined in a 65-nm CMOS development. The proposed gatherer showed the abatement of locale and the force use by 17% and 19%, exclusively, stood out and the authority from the conventional CLA snake based structure. By virtue of the MAC structure, the proposed arrangement decreased both the district and force by 20%. We acknowledge that the proposed idea to utilize the stand-out typical for AI count for logically powerful MAC design can be gotten in various neural framework hardware reviving operator structures later on.

## REFERENCES

[1] A. Krizhevsky, I. Sutskever, and G. E. Hinton, "Image Net classification with deep convolution neural networks," in Proc. Adv. Neural Inf. Process. Syst., 2012, pp. 1097–1105.

[2] K. Simonyan and A. Zisserman. (2014). "Very deep convolution networks for large-scale image recognition." [Online]. Available: https://arxiv.org/abs/1409.1556

[3] C. Szegedy et al., "Going deeper with convolutions," in Proc. IEEE Conf. Comput. Vis. Pattern Recognit., Jun. 2015, pp. 1–9.

[4] A. Graves, A.-R. Mohamed, and G. Hinton, "Speech recognition with deep recurrent neural networks," in Proc. IEEE Int. Conf. Acoust., Speech Signal Process. (ICASSP), May 2013, pp. 6645–6649.

[5] J. K. Chorowski, D. Bahdanau, D. Serdyuk, K. Cho, and Y. Bengio, "Attention-based models for speech recognition," in Proc. Adv. Neural Inf. Process. Syst., 2015, pp. 577–585.

[6] Y. H. Chen, T. Krishna, J. S. Emer, and V. Sze, "Eyeriss: An energy-efficient reconfigurable accelerator for deep convolution neural networks," IEEE J. Solid-State Circuits, vol. 52, no. 1, pp. 127–138, Jan. 2017.

[7] B. Moons, R. Uytterhoeven, W. Dehaene, and M. Verhelst, "Envision: A 0.26-to-10tops/w sub word-parallel dynamic-voltage-accuracy frequency-scalable convolution neural network processor in 28nm FDSOI," in Proc. IEEE Int. Solid-State Circuits Conf. (ISSCC), Feb. 2017, pp. 246–247

[8] C. S. Wallace, "A suggestion for a fast multiplier," IEEE Trans. Electron. Comput., vol. EC-13, no. 1, pp. 14–17, Feb. 1964.

[9] L. Dadda, "Some schemes for parallel multipliers," Alta Frequenza, vol. 34, no. 5, pp. 349–356, Mar. 1965.

[10] P. F. Stelling and V. G. Oklobdzija, "Implementing multiply-accumulate operation in multiplication time," in Proc. 13th IEEE Symp. Comput. Arithmetic, Jul. 1997, pp. 99–106.

[11] K. K. Parhi, VLSI Digital Signal Processing Systems: Design and Implementation. New Delhi, India: Wiley, 1999.

[12] T. T. Hoang, M. Sjalander, and P. Larsson-Edefors, "A high-speed, energy-efficient two-cycle multiply-accumulate (MAC) architecture and its application to a double-throughput MAC unit," IEEE Trans. Circuits Syst. I, Reg. Papers, vol. 57, no. 12, pp. 3073–3081, Dec. 2010.

[13] W. J. Townsend, E. E. Swartzlander, and J. A. Abraham, "A comparison of Dadda and Wallace multiplier delays," Proc. SPIE, Adv. Signal Process. Algorithms, Archit., Implement. XIII, vol. 5205, pp. 552–560, Dec. 2003, doi: 10.1117/12.507012.

[14] M. Courbariaux, Y. Bengio, and J.-P. David, "Binary Connect: Training deep neural networks with binary weights during propagations," in Proc. Adv. Neural Inf. Process. Syst., 2015, pp. 3123–3131.

[15] M. Rastegari, V. Ordonez, J. Redmon, and A. Farhadi, "XNOR-Net: Image Net classification using binary convolution neural networks," in Proc. Eur. Conf. Comput. Vis. Springer, 2016, pp. 525–542.

[16] M. Gao, J. Pu, X. Yang, M. Horowitz, and C. Kozyrakis, "Tetris: Scalable and efficient neural network acceleration with 3d memory," in Proc. 22nd Int. Conf. Archit. Support Program. Lang. Oper. Syst., 2017, pp. 751–764.

[17] A. Parashar et al., "SCNN: An accelerator for compressed-sparse convolution neural networks," in Proc. 44th Annu. Int. Symp. Comput. Archit., Jun. 2017