# Implementation of DAC for Energy Efficient SRAM with In-memory Computation

**Proposal**

**Prepared By:**

**Kazi Barria Nine**
Associate Engineer
Neural Semiconductor Ltd.

Date: 30th August 2021

NS
NEURAL SEMICONDUCTOR

# Contents

# 1 Abstract:

This proposal states a particular idea for implementing a Digital to Analog Converter (DAC) which has application in column wise conversion of digital codes for SRAM with in-memory Multiply and Accumulate (MAC) computation (1). The design features the use of Delay Locked Loop (DLL) that generates well defined delayed signals that are used to generate timing signals used by a Time to Analog Circuit to produce the analog output. The architecture offers the opportunity to scale the no. of channels with minimal increase in area as the delayed signals are global and shared with every channel.

# 2 Background:

This DAC will be used in energy efficient convolutional SRAM with in-memory computing. Detailed info can be found in reference 1, however a brief summary is presented here for convenience. Stated CSRAM array consists of $256\ rows \times 64\ columns$ of SRAM bit-cell. It is divided into 16 local arrays with 16 rows each. The Input digital values ($X_{IN}$) are processed by columnwise DACs, which convert the digital $X_{IN}$ codes to analog input voltages on the global read bit lines (GRBLs). The GRBLs are shared by all the local arrays, implementing the fact that in CNNs each input is shared/processed in parallel by multiple filters. The DLL circuit is used for generating pulse, it is shared by all 64 DACs. During the first phase of the SRAM operation the digital convolution input ($X_{IN}$) is converted into an analog voltage ($V_a$) using a column wise DAC. The analog voltage is used to pre-charge the global read bit line (Analog Output in proposed block diagram). This bitline will be shared by all 16 local Arrays to perform the addition part of the MAC operation.
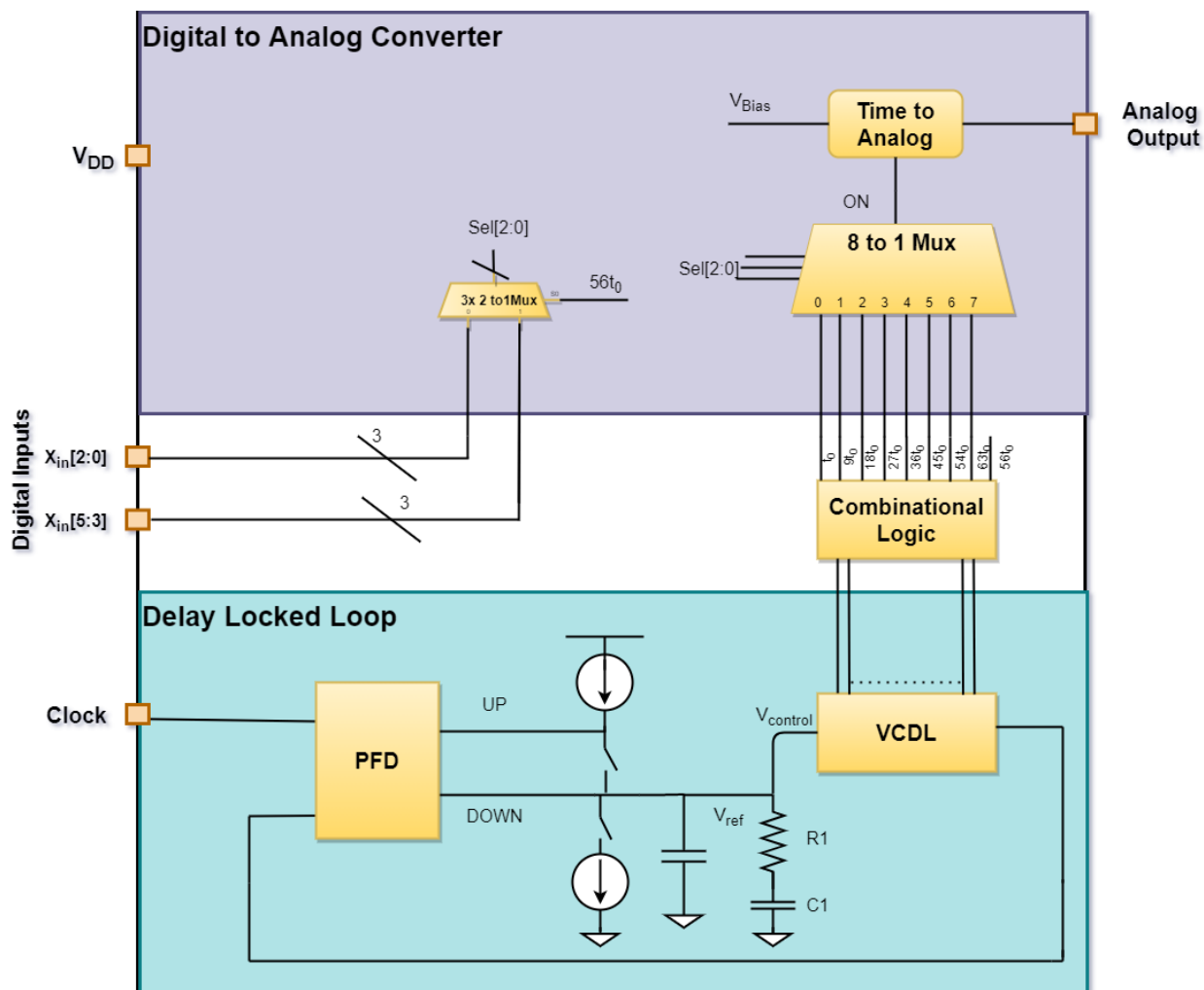
# 3 Proposed Block Diagram:



*Figure 1 Block Diagram of Proposed Arcchitecture*

# 4   Circuit Description:

The proposed block diagram is stated above which includes three sub-blocks – DLL (Delay Locked Loop), Combinational logic block (for converting DLL outputs in appropriate global timing signals) and lastly, the actual time to analog part to convert digital input to analog voltage. Fig. 1 depicts the block diagram for the proposed architecture.

## 4.1   Time to Analog Conversion:

A cascode pMOS stack biased in saturation region is used as a current source. The capacitor on the Analog output line is charged with this fixed current for a specific time which is determined by the ON pulse width which in turn depends on the input code. This ON pulse width is created from digital input code ($X_{IN}[5:0]$) and several control signals ($t_0$ to $63t_0$) using a digital-to-time converter implemented using a combination of DLL & MUXes which multiplexes the delayed signals according to the input code. A single continuous ON pulse is required for achieving good linear relationship between $V_a$ vs $X_{IN}$.
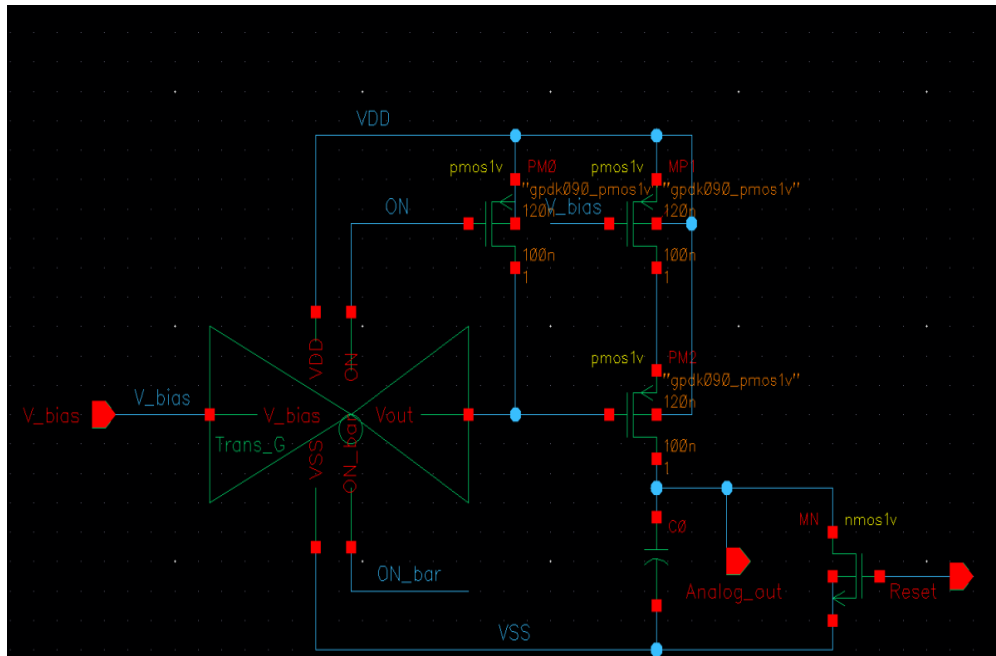


*Figure 2 Time to Analog Circuit*

## 4.2   Multiplexers:

A conventional multiplexer will be used to generate the ON signal required by the Time to Analog converter circuit, selecting between generated delay signals from the DLL and Combinational Circuits depending on
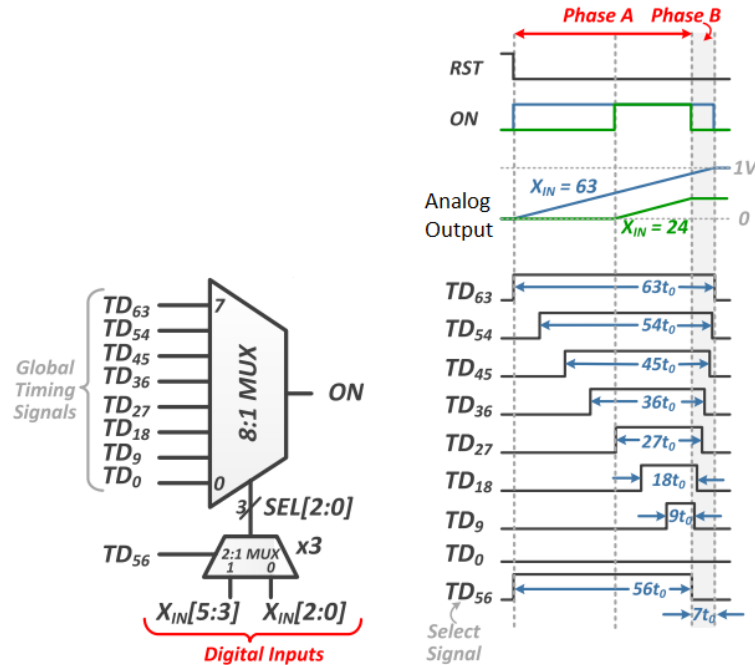


Figure 3 MUX Architecture [Taken from Reference 1]

## 4.3   Combinational Circuit:

A series of NOR & NAND gates are to be used to generate the required signals from a sequence of delayed signals from a DLL.

## 4.4    Delay Locked Loop:

A DLL is used to generate a sequence of square waves each delayed from the last by an amount defined by the time required to charge the output capacitor by 1LSB. The DLL is a closed loop circuit consisting of a Phase Frequency Detector (PFD), a Charge Pump & Filter Circuit and a Voltage Controlled Delay Line (VCDL). The use DLL ensures that the delay is well regulated and resilient against PVT variations and minimal effects of jitter.

### Phase Frequency Detector:

The PFD consists of two D flip flops with asynchronous reset and an AND gate. Its job is to compare the output of the VCDL (feedback signal) with a reference clock, and generate UP and DOWN pulses to drive the downstream Charge Pump, depending on whether the feedback signal leads or lags the reference clock.
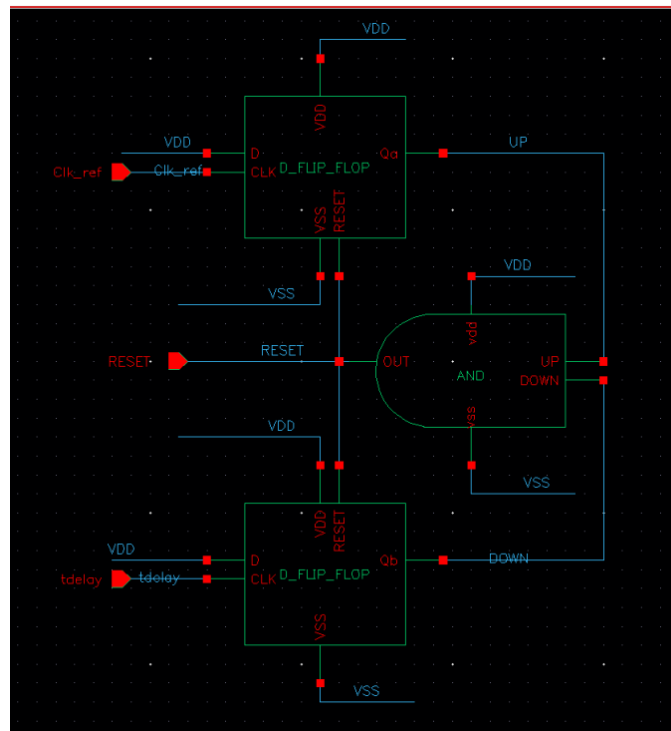


*Figure 4 Phase Frequency Detector*

### 4.4.1 Charge Pump & Loop Filter:

This block generates the appropriate control voltage for the VCDL to provide the required delay regulated by the DLL loop. It is done by injecting or removing charge from the filter capacitor by means of two current sources implemented with FETs, according to the UP & DOWN signals from the PFD. The diagram shows ideal sources as reference, these are to be replaced by reference current circuits.
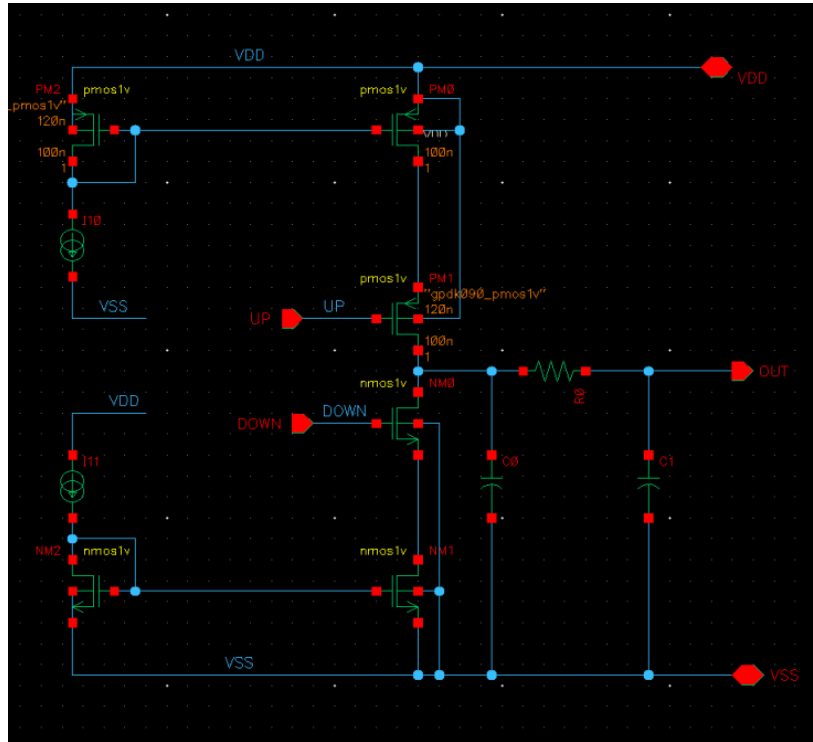


*Figure 5 Charge Pump and FIlter*

### 4.4.2 Voltage Controlled Delay Line:

This block generates a sequence of delayed square waves where the delay is proportional to the input control voltage. This is implemented using several current starved inverter stages where the current is tuned according to the input voltage that tunes the delay.
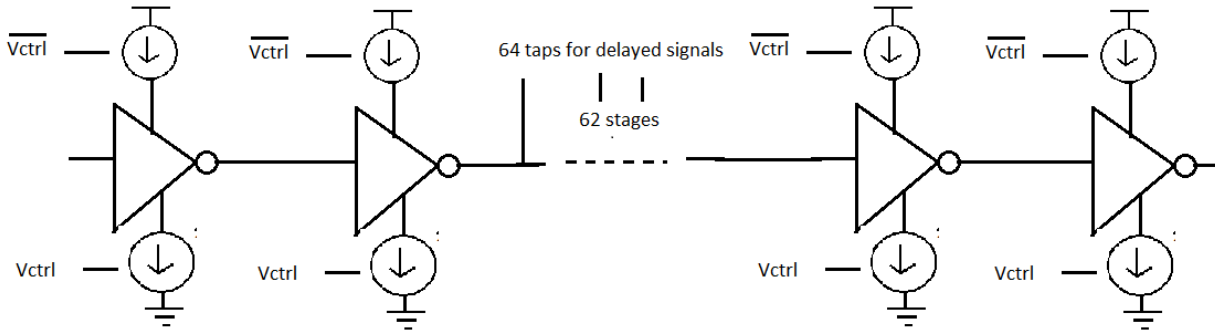


*Figure 6 Voltage Controlled Delay Line*

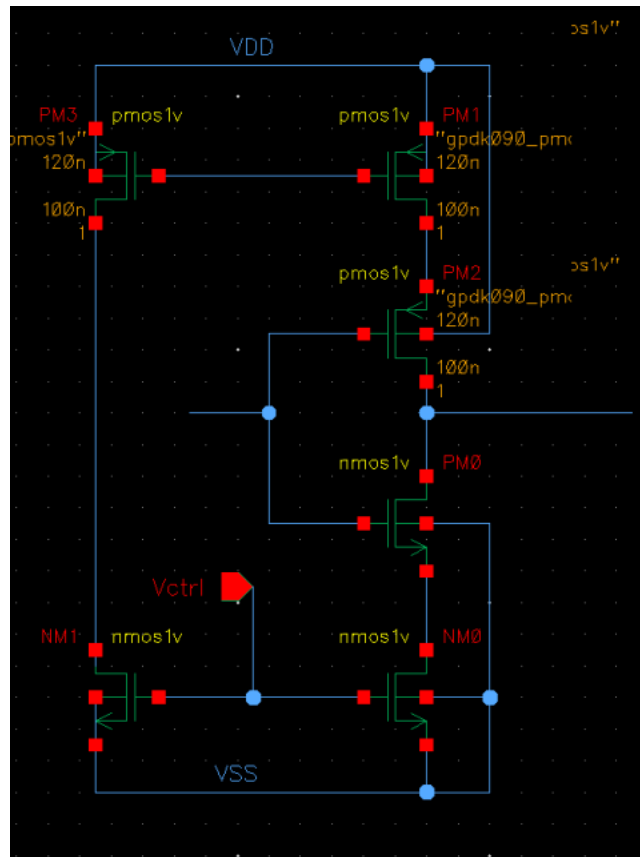This circuit for each inverter stage is given below:



*Figure 7 VCDL Current Starved Inverter Stage*

## 5 Expected Performance:

| Parameter | Value | Units |
|---|---|---|
| Supply | 1.8 | V |
| Full Scale voltage | 800 | mV |
| Number of bits | 6 | bits |
| 1LSB Voltage (V1LSB) | 12.7 | mV |
| Conversion Time | 100 | ns |
| Frequency | 10 | MHz |
| INL | <2 | LSB |
| DNL | <2 | LSB |
| Offset | <1 | LSB |
| Gain Error | <1 | % |

# 6 Acknowledgement & Reference:

1. Biswas and A. P. Chandrakasan, "CONV-SRAM: An Energy-Efficient SRAM With In-Memory Dot-Product Computation for Low-Power Convolutional Neural Networks," in *IEEE Journal of Solid-State Circuits*, vol. 54, no. 1, pp. 217-230, Jan. 2019, doi: 10.1109/JSSC.2018.2880918.

# 7 Team Members:

- Bustana Teene
- Tasnim Hasan
- Hasin Akhteyar
- Mayeesha Fairuz Ahmed
- Rouhan Noor
- Saleh Omar
- Nugaira Gahan