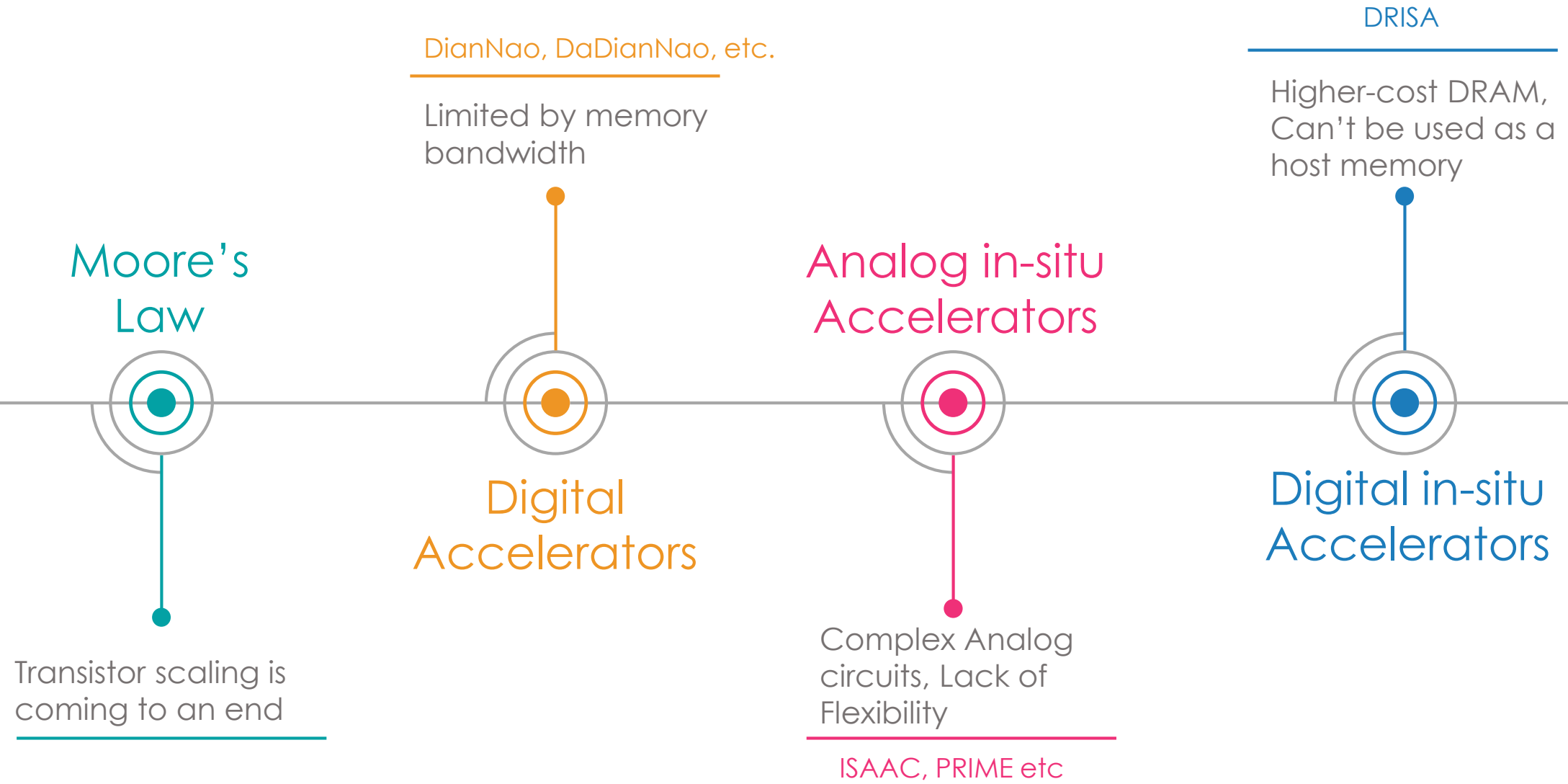


Moving CNN Accelerator Computations Closer to Data



Sumanth Gudaparthi
Surya Narayanan
Rajeev Balasubramonian

Evolution of CNN Accelerators



SRAM based In-Situ Computation Accelerator



DA vs SISCA

Perform
Computations In-
Situ



AIA vs SISCA

Use SRAM cells
to perform In-
Situ
Computations

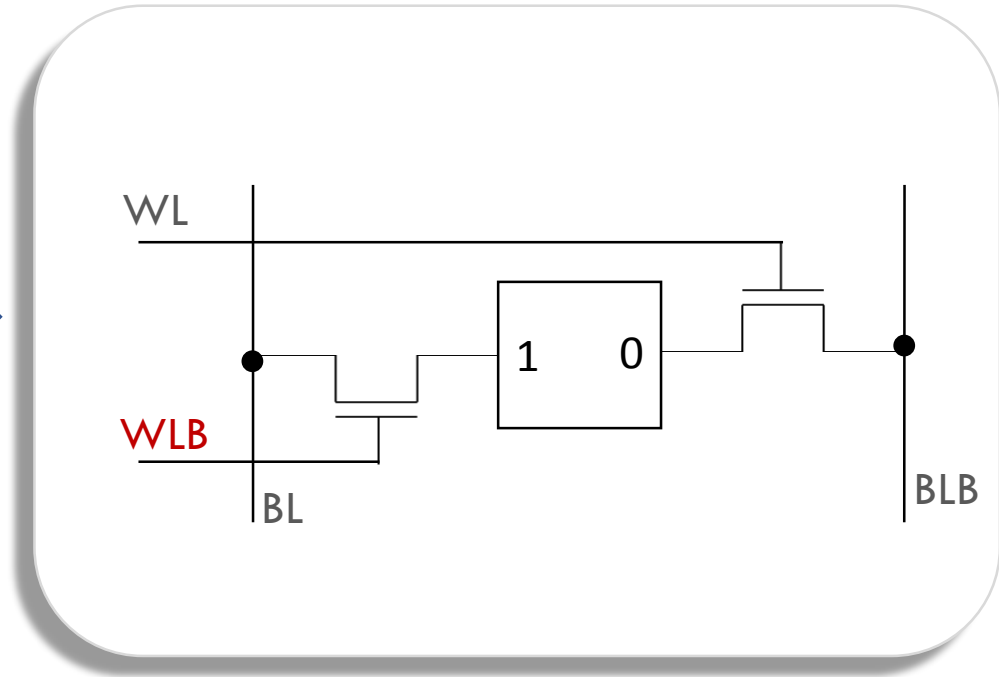
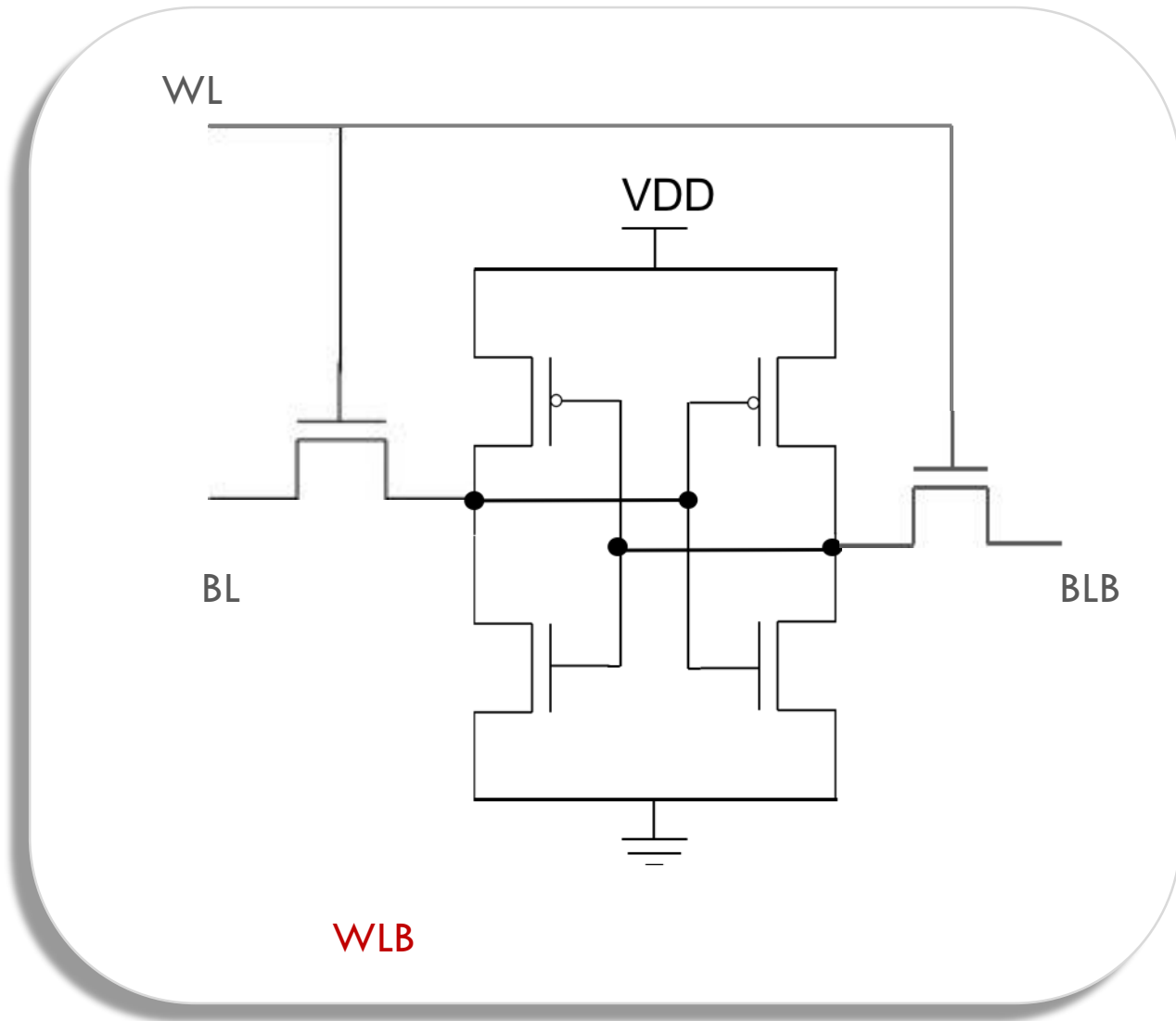


DIA vs SISCA

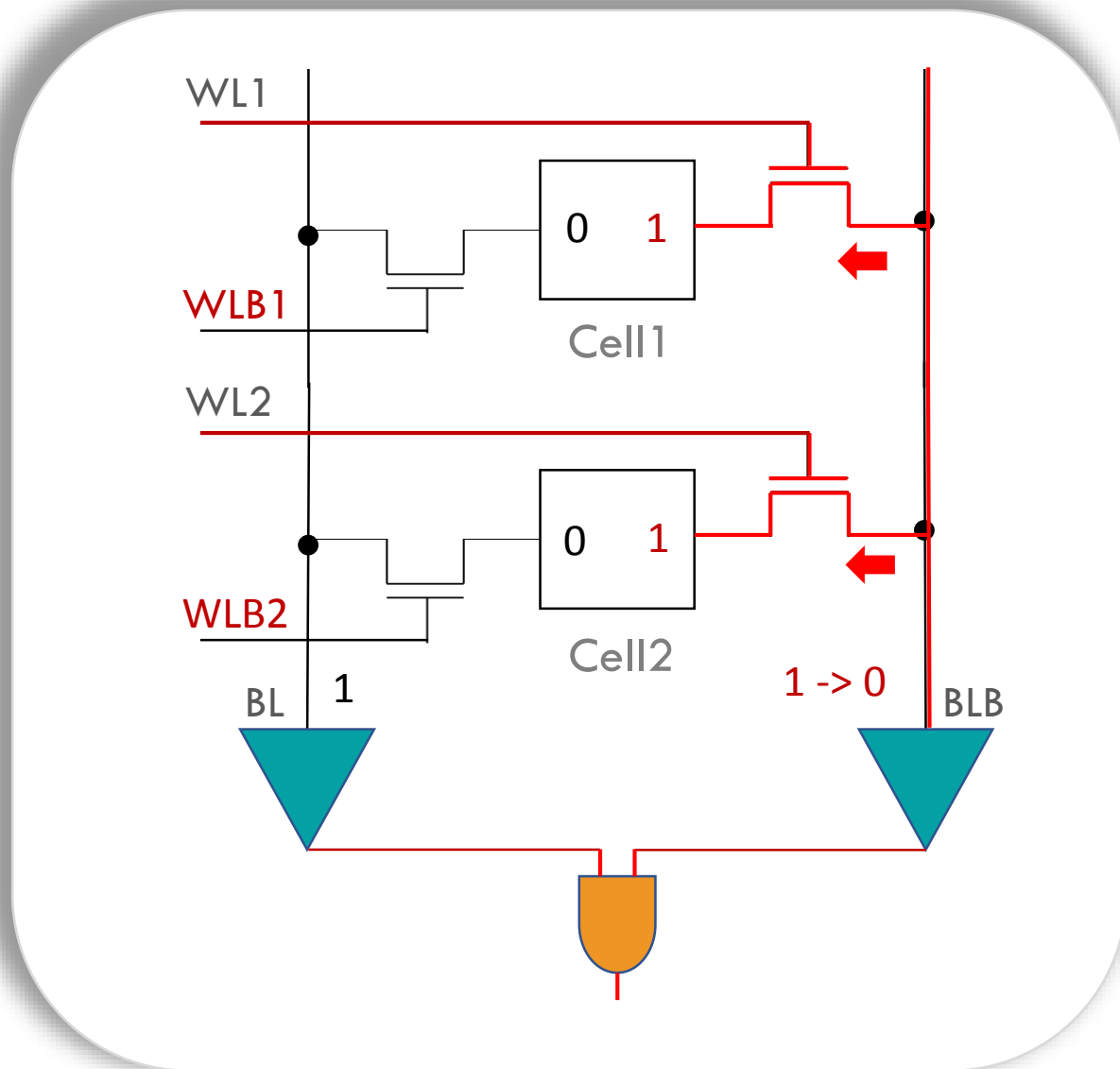
Modify the LLC.
Trivial overhead
on baseline Cache
operations



Logic-In-Memory



Logic-In-Memory



- Pre-charge the bit-lines
- Activate the word-lines
- Discharge of bit-line voltage through Cell1
- Discharge of bit-line voltage through both Cells
- Bit-line stays Pre-charged

Enabling In-Situ Multiplication in Caches

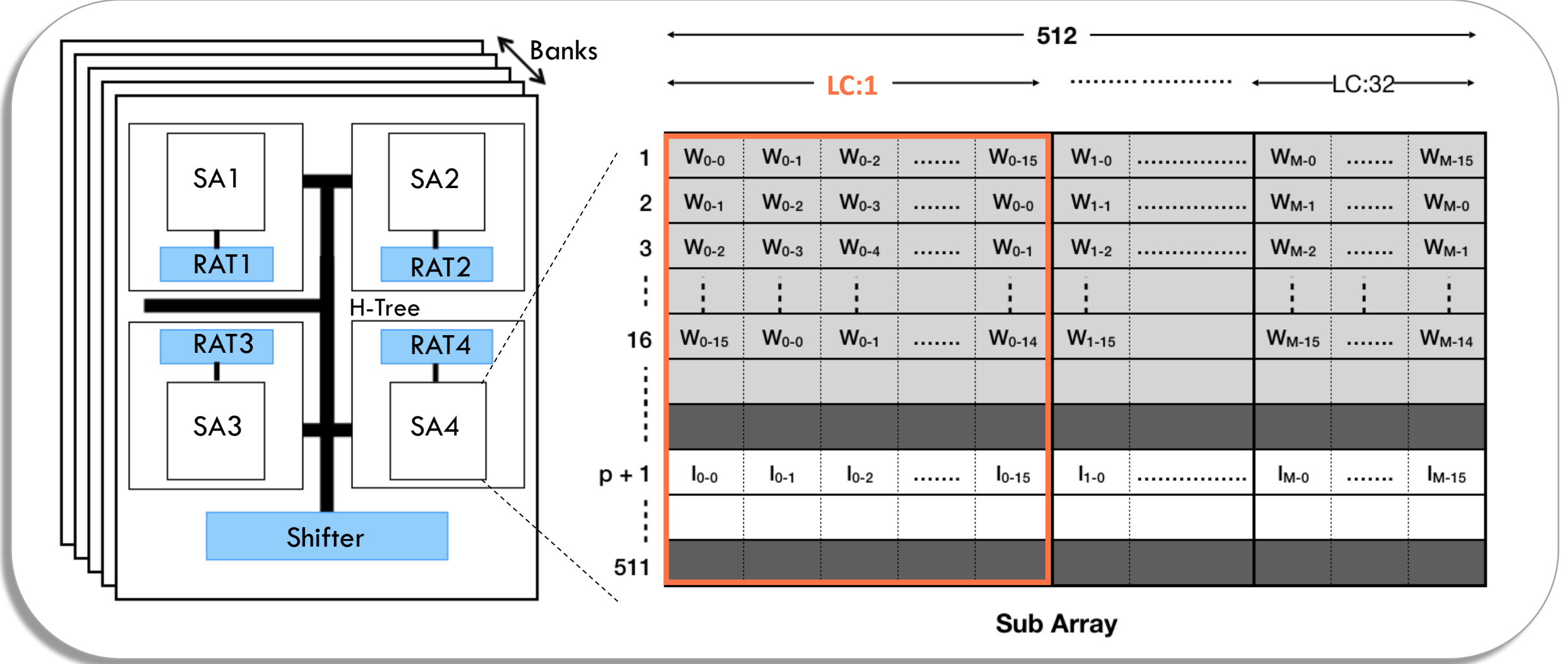


$$\begin{array}{r}
 W_{0-0} \quad W_{0-1} \quad W_{0-2} \\
 * \quad I_{0-0} \quad I_{0-1} \quad I_{0-2} \\
 \hline
 \\
 \\
 W_{0-0}I_{0-2} \quad W_{0-1}I_{0-2} \quad W_{0-2}I_{0-2} \\
 \\
 W_{0-0}I_{0-1} \quad W_{0-1}I_{0-1} \quad W_{0-2}I_{0-1} \\
 \\
 W_{0-0}I_{0-0} \quad W_{0-1}I_{0-0} \quad W_{0-2}I_{0-0}
 \end{array}$$

W_{0-0}	W_{0-1}	W_{0-2}
W_{0-1}	W_{0-2}	W_{0-0}
W_{0-2}	W_{0-0}	W_{0-1}
I_{0-0}	I_{0-1}	I_{0-2}

C_{a-b} : bit number-b in a^{th} variable of C

SISCA Organization

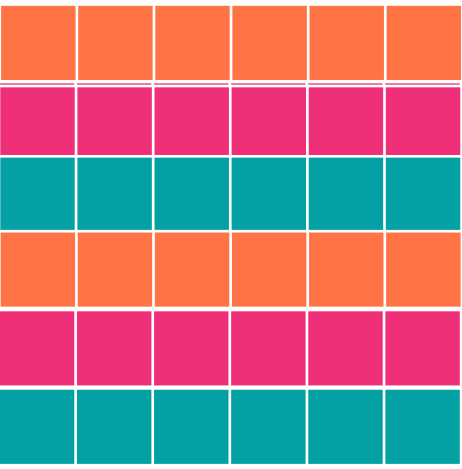


C_{a-b} : bit number-b in a^{th} variable of C

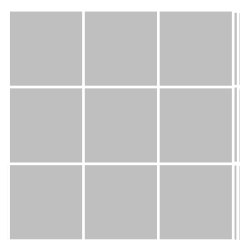
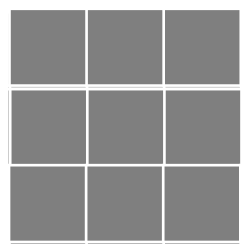
Legend:

- Unused Sub-array Entries
- Kernel Entries
- Feature Map Entries

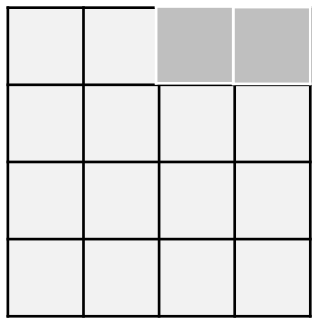
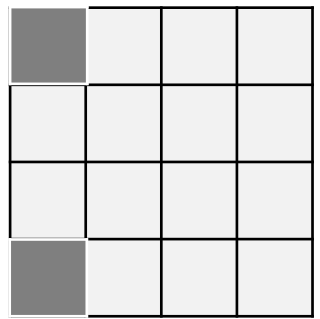
SISCA Dataflow



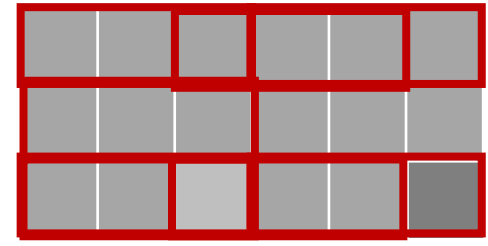
Input Feature Map
(6x6)



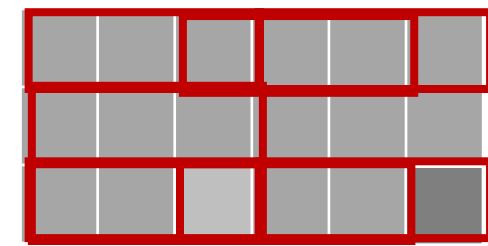
Kernel Maps
2x(3x3)



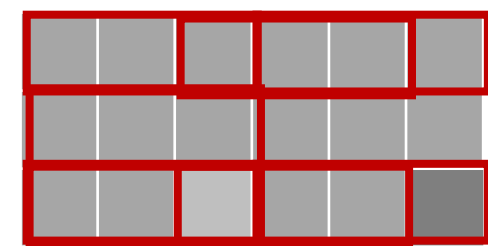
Output Feature Maps
2x(4x4)



Sub Array 1

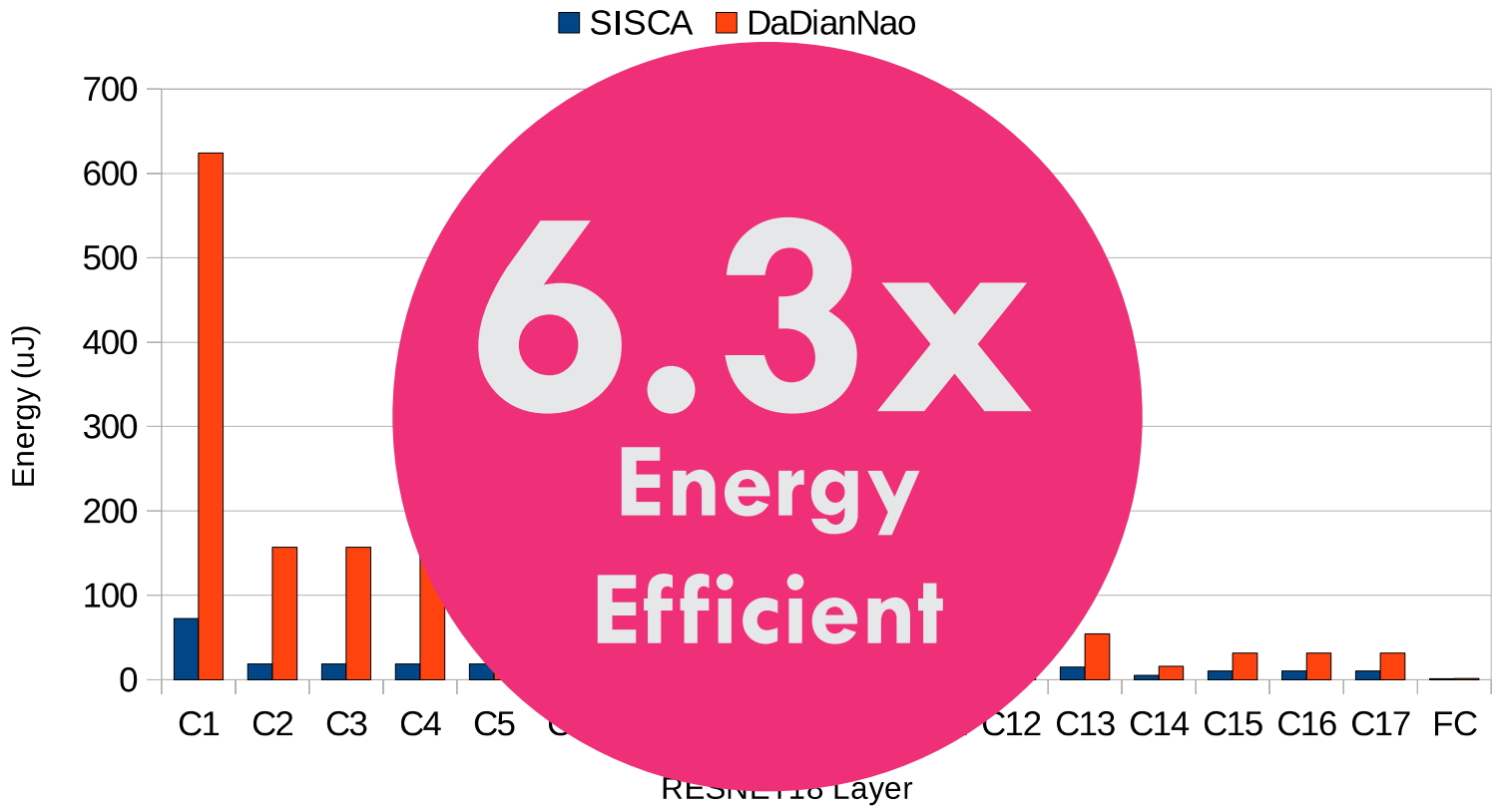


Sub Array 2

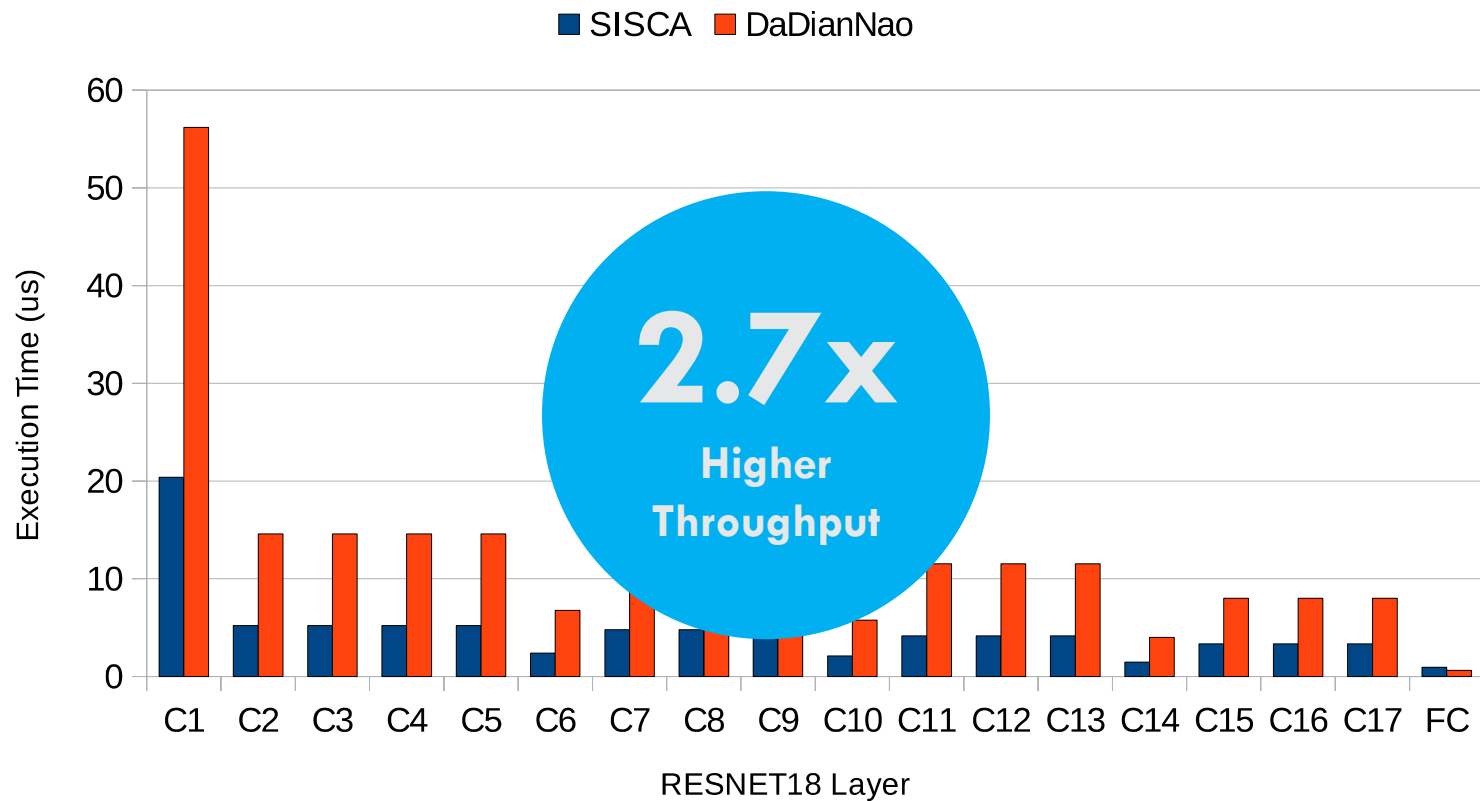


Sub Array 3

Energy Improvements



Performance Improvements



Conclusions and Future Work



- SISCA is an SRAM in-situ computation Engine for Convolution Neural Networks
- Uses on-chip Last Level Cache (LLC) to perform computations
- SISCA is 6.3x Energy efficient, and has 2.7x higher throughput than DaDianNao
- Better dataflow and mapping mechanisms can further improve the Energy and Throughput.
- Need to work on better scheduling mechanisms to distribute the general purpose workload, and CNN data across the Cache.

Questions?