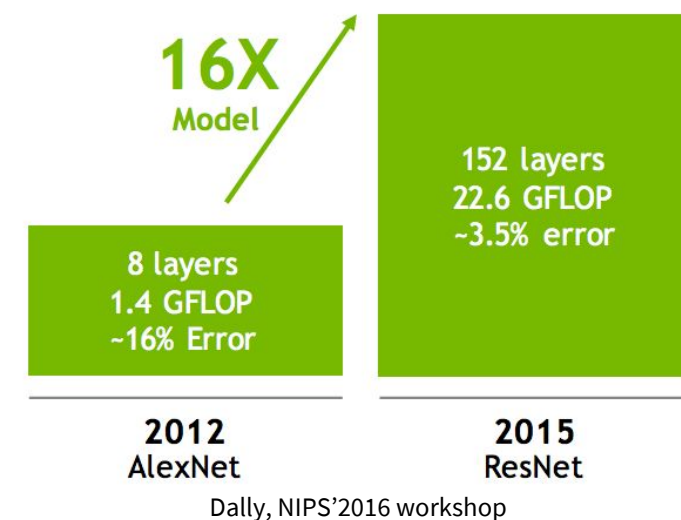
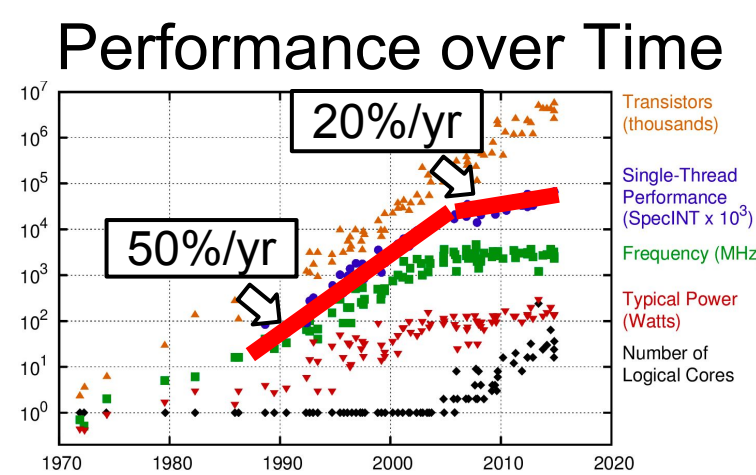


## Background

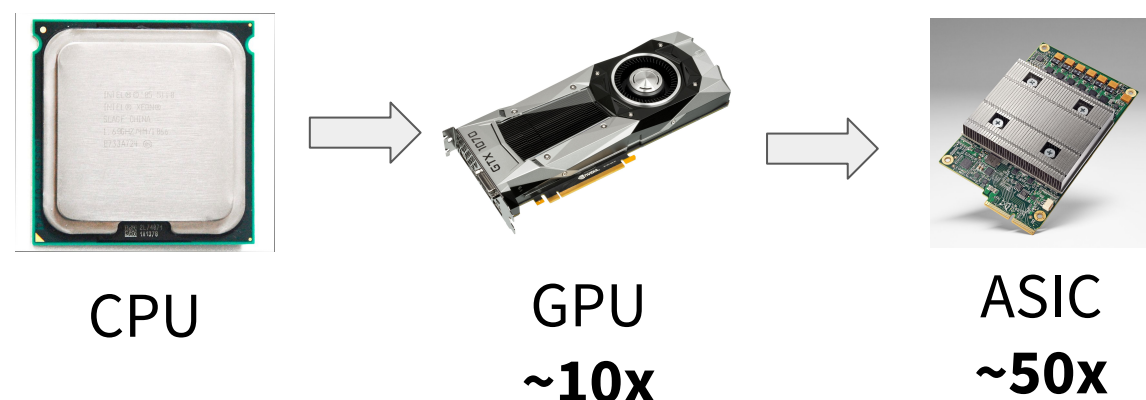
- New CNN models require massive computation



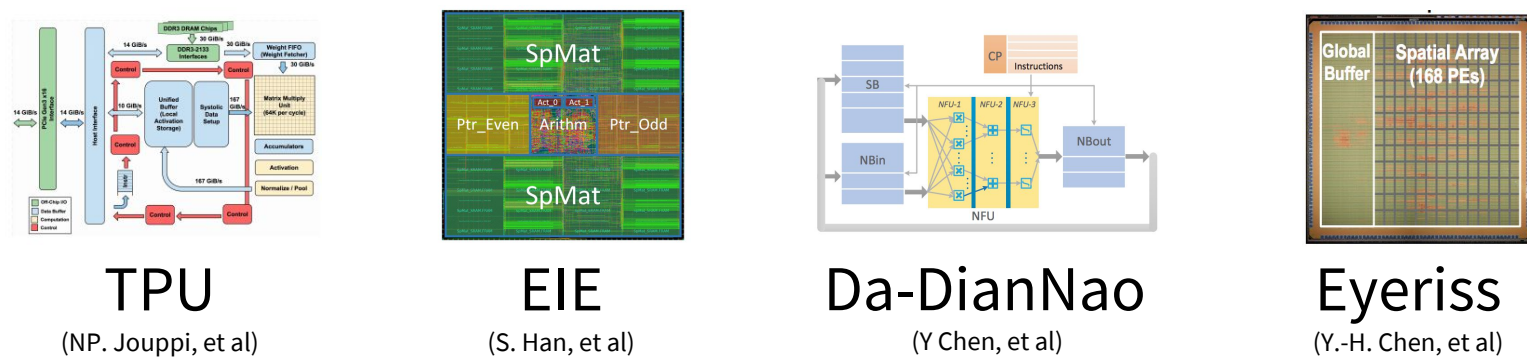
- Hardware Improvements have slowed



- One solution: specialized hardware



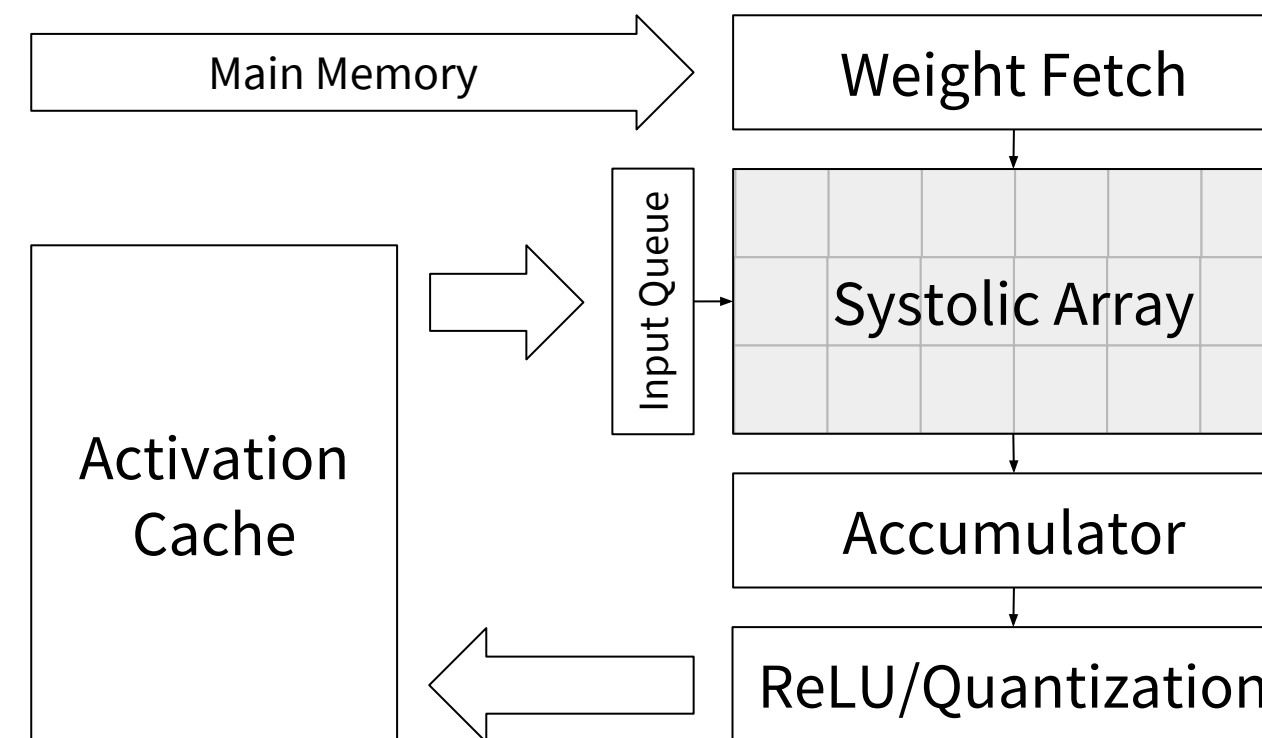
- Prior Work



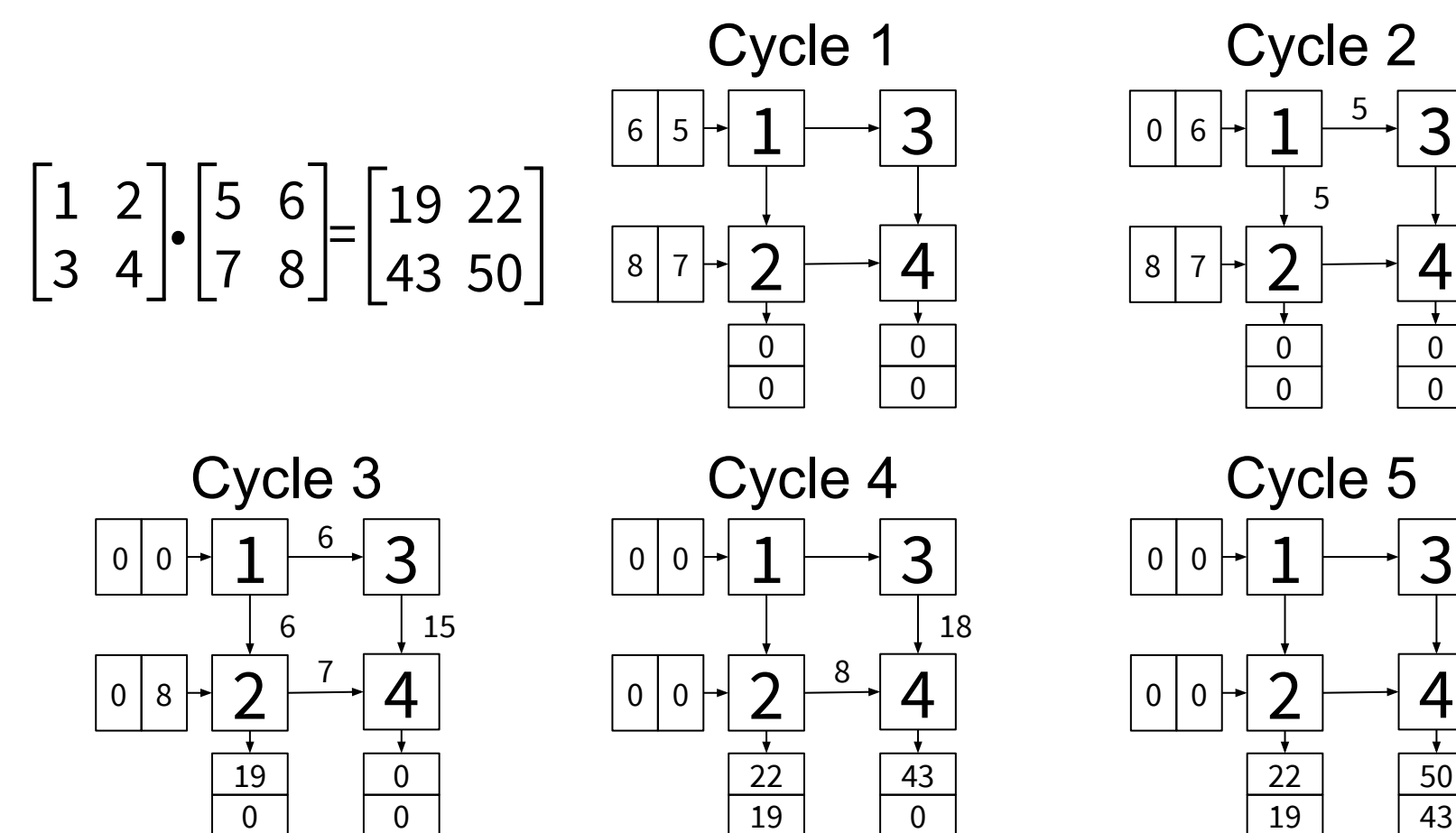
## Problem Statement

- Design and implement a hardware CNN accelerator
- Compare our design's power and throughput to TPU and a GPU

## Design

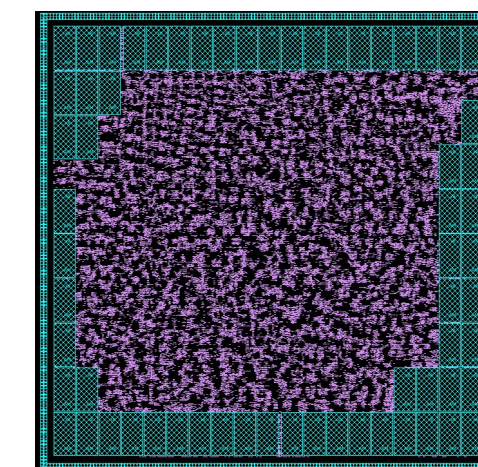


- Main component is systolic array matrix multiply
  - Conv, Batch Norm, FC, etc, can be transformed to dense matrix mult.
- Activations cached for fast reuse



- Weight stationery design
  - Weights loaded from memory before multiply
  - Multiple inputs can reuse weights without reloading

## Evaluation



Final layout PnR in 32nm

- Evaluation on 8-bit quantized SqueezeNet (78.4 Top-5 ImageNet)

	Power (W peak)	Clock (MHz)	TOP/s		Area (mm <sup>2</sup> )
			Int	FP	
K80/die	98	560	-	2.8	561
TPU/die	40	700	92	-	<300
Ours	<b>20.9</b>	<b>714</b>	<b>93.6</b>	-	<b>76.6</b>

- Power breakdown

Memory	3.2W	15.7%
Clock	1.1W	5.4%
Register	7.3W	35%
Combinational	9.1W	43%

## Conclusion

- Our design efficiently accelerates CNN inference
  - 33x TOPS @ 0.21x power vs GPU
- Future: accelerate sparse matrix ops