

# Approximate Computing and Microfluidic Cooling for Enhanced Machine Learning

**Hardik Sharma**

†**William Wahby**

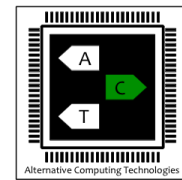
†Thomas Sarvey

†Muhannad S. Bakir

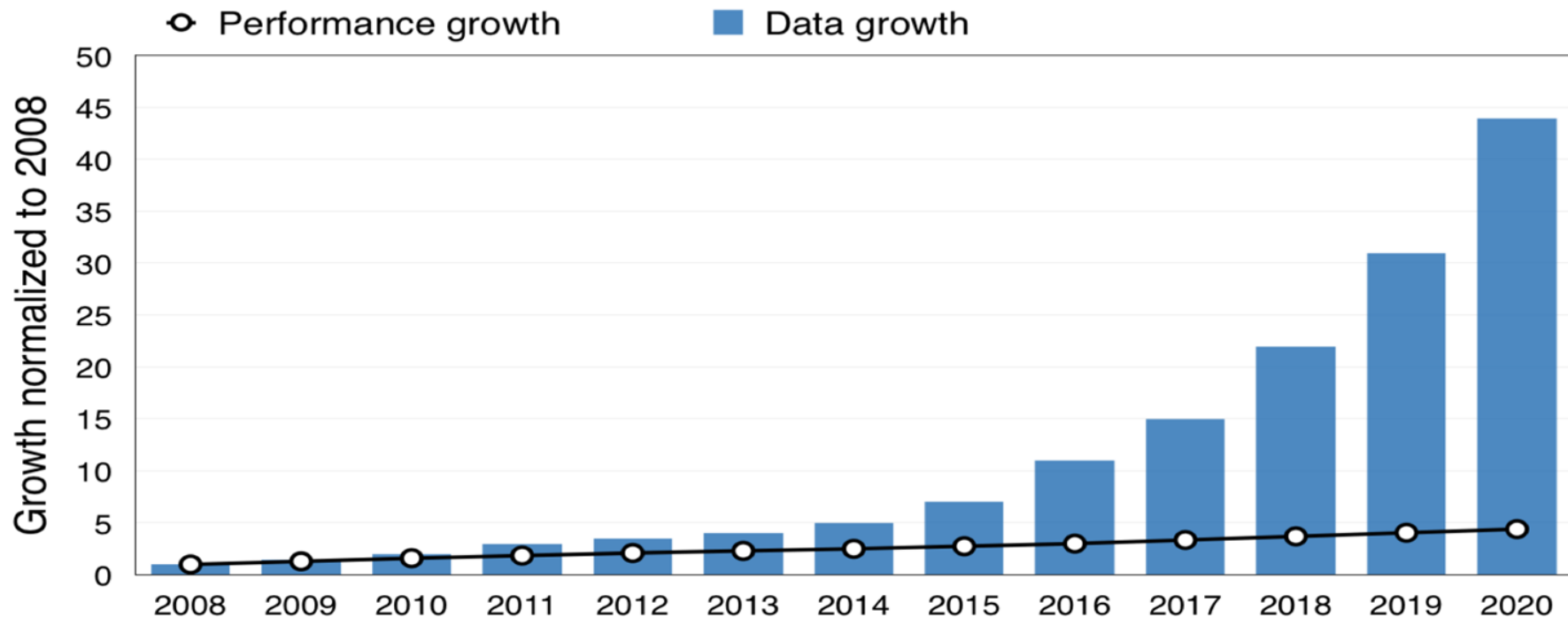
Hadi Esmaeilzadeh

Alternative Computing Technologies (**ACT**) Lab  
Georgia Institute of Technology

†Integrated 3D Systems (**I3DS**) Group  
Georgia Institute of Technology

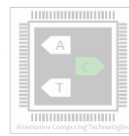


# Data growth vs. Performance

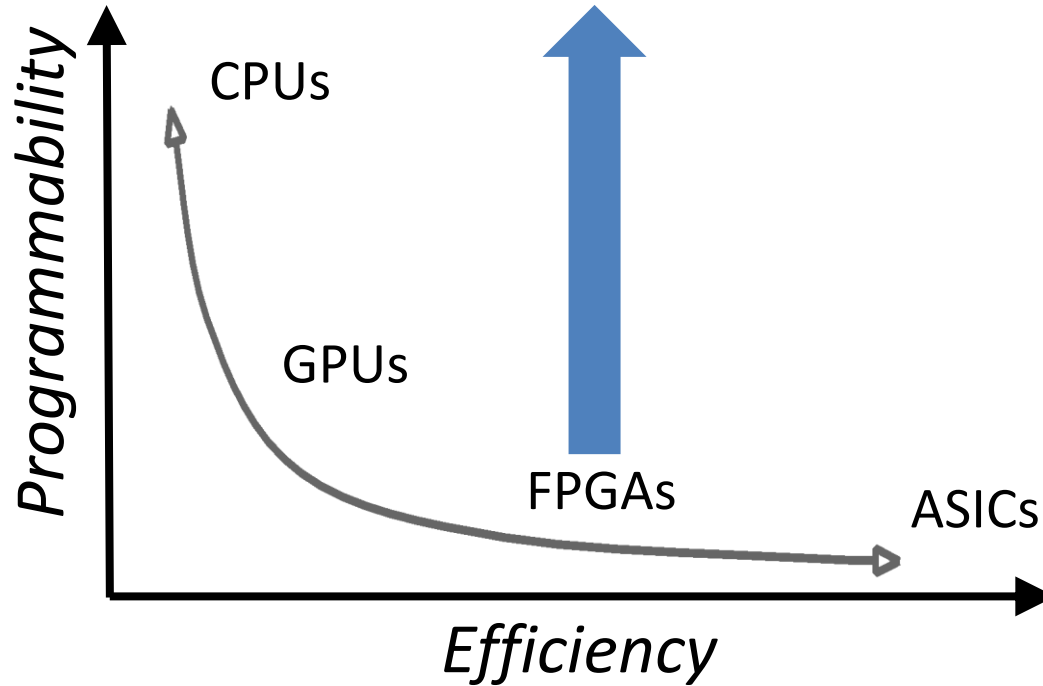


Data growth trends: IDC's Digital Universe Study, December 2012

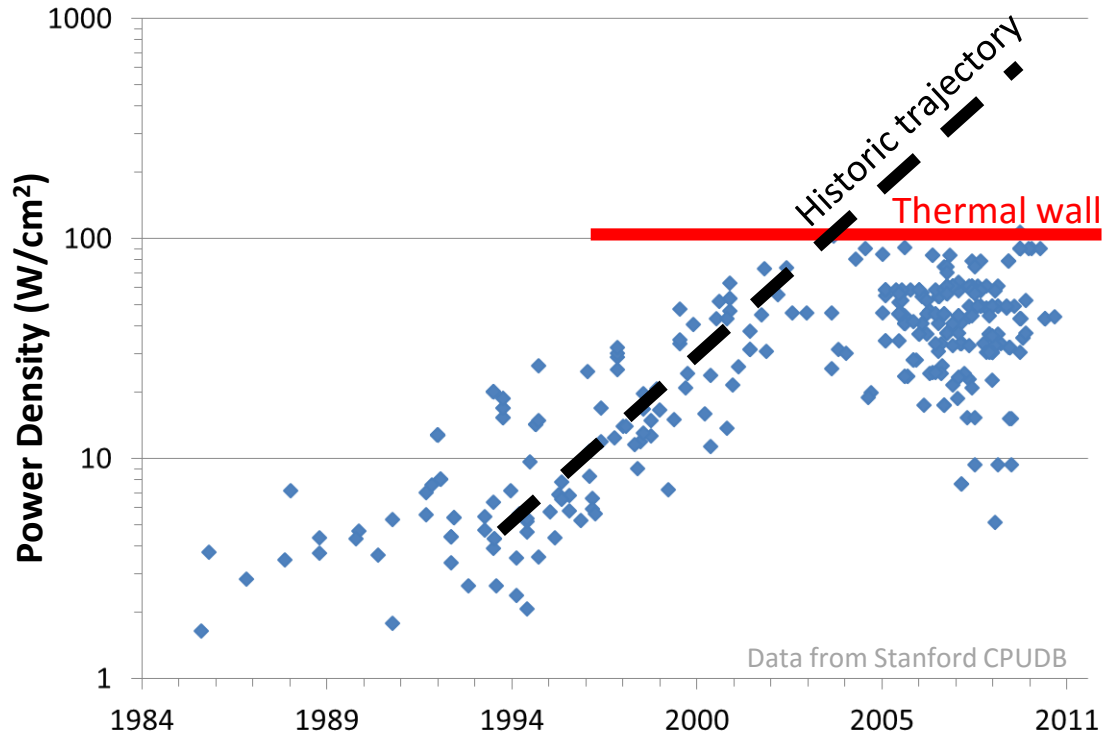
Performance growth trends: Esmailzadeh et al, "Dark Silicon and the End of Multicore Scaling," ISCA 2011



# Programmability vs. Efficiency

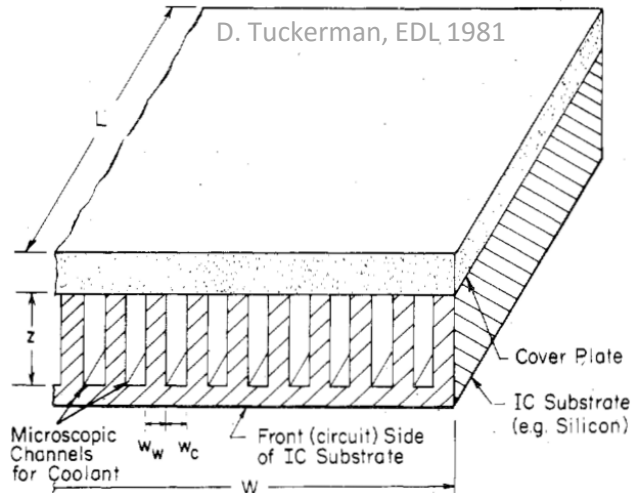


# Heat limits system performance

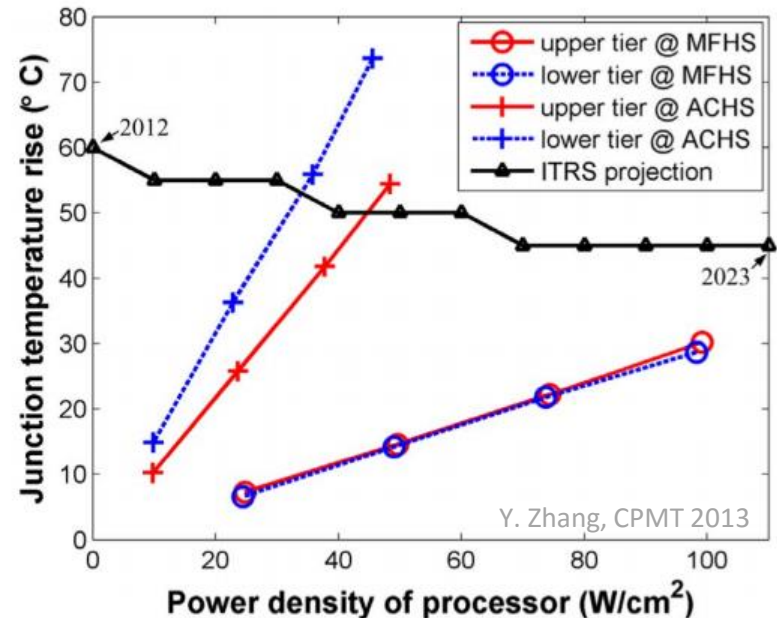


# Better cooling improves performance

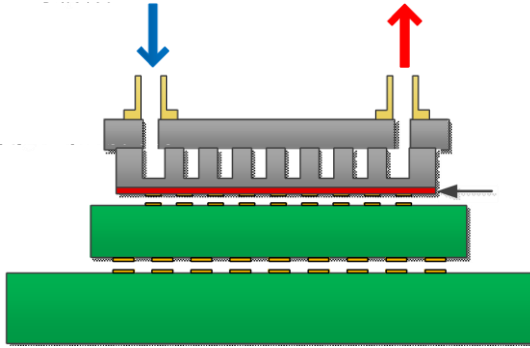
First Proposed in 1981



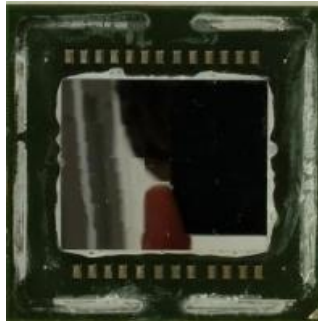
Non-functional characterization in 2013



# Microfluidically-cooled Stratix V FPGA

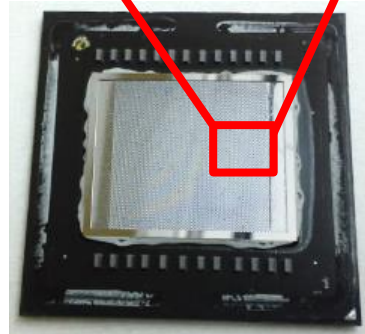
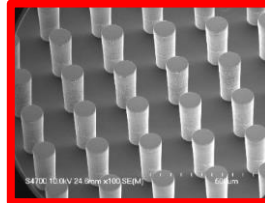


Original Die



Delidded die

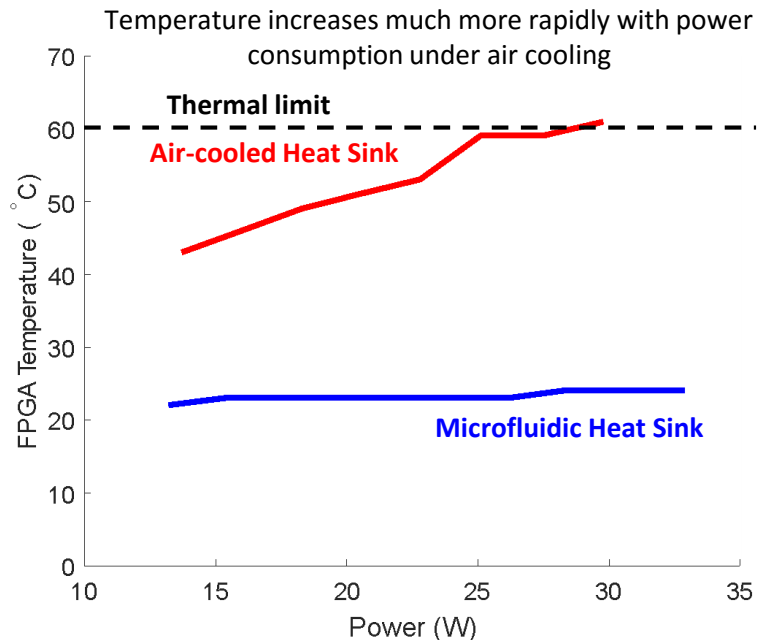
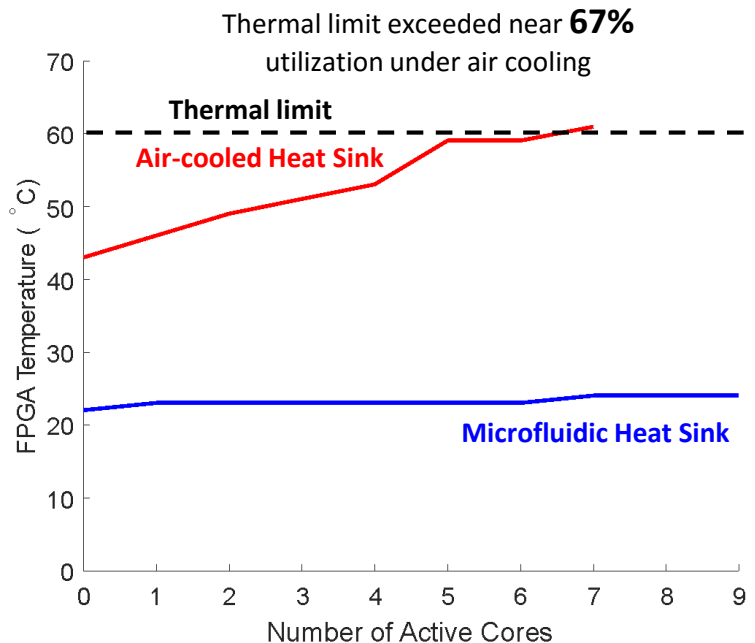
Micropinfin closeup



Etched Microfluidic Heat Sink



# Microfluidically-cooled Stratix V FPGA



Junction-to-ambient  $R_{th} \approx 0.08^{\circ}\text{C}/\text{W}$

Nominally expect **only 40 °C increase** over ambient at **500W** power dissipation

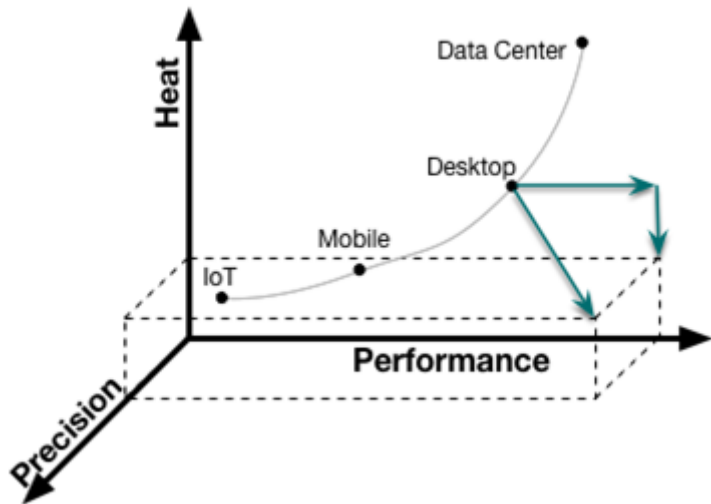
# Approximate Computing for FPGA Acceleration

Deep Neural Networks have high tolerance to approximation.

(DeepCompression ICLR2016)

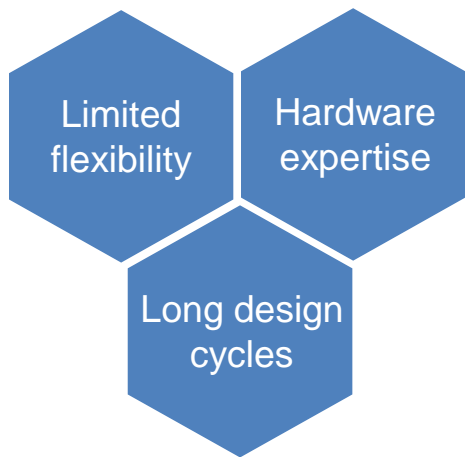
**Relaxing Precision will yield higher performance.**

- Increased parallelism through reduced resource usage.
- Reduce bandwidth.

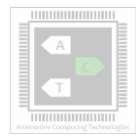
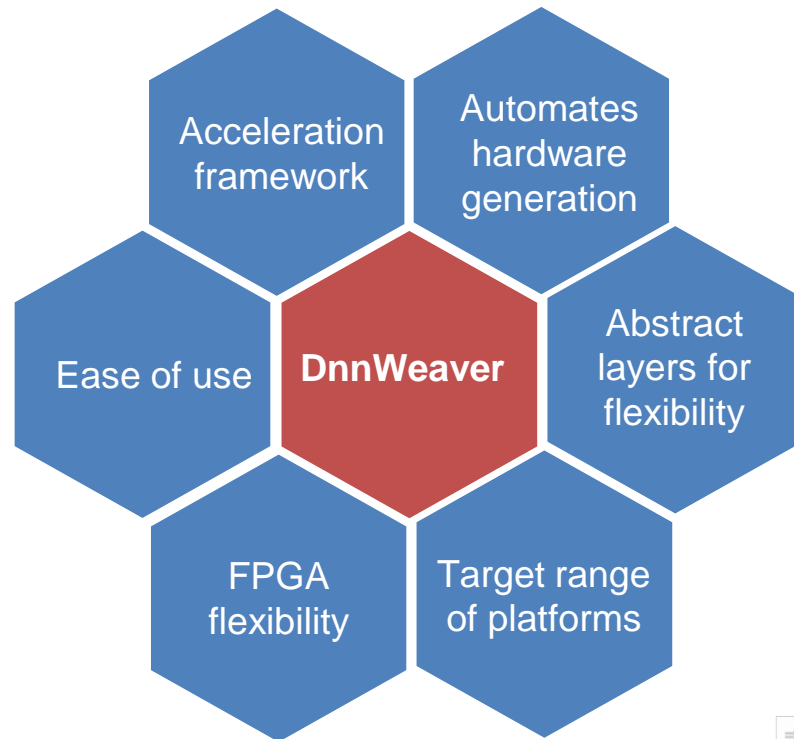




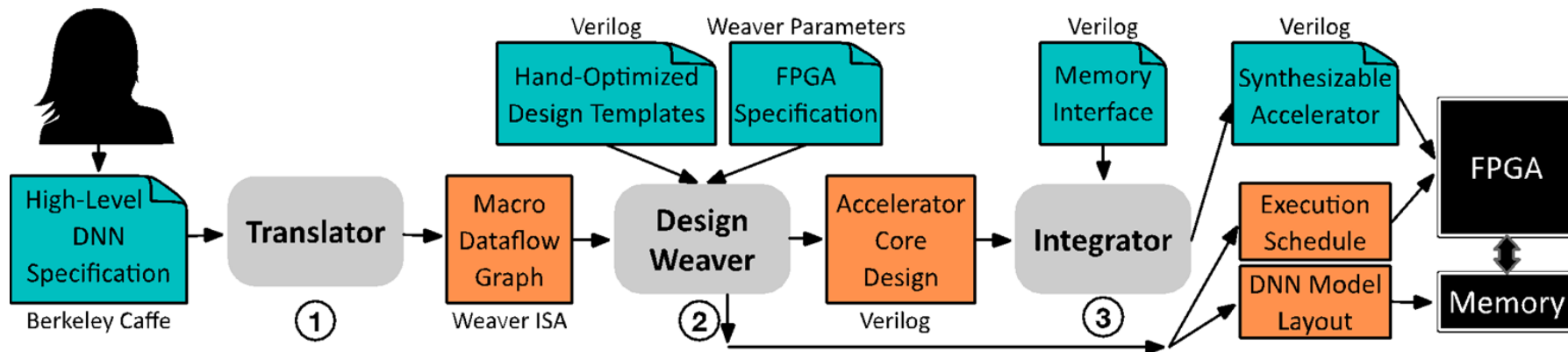
# Challenges in Hardware Acceleration



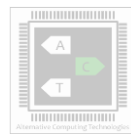
DnnWeaver



# Compilation flow

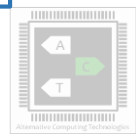


We present a comprehensive framework for accelerating DNNs from high-level abstractions



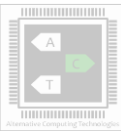
# Benchmarks

<b>LENET</b>	<b>Character recognition</b>	<b>C-&gt;P-&gt;C-&gt;P-&gt;I-&gt;A-&gt;I</b>
<b>Siamese</b>	<b>Character recognition</b>	<b>C-&gt;P-&gt;C-&gt;P-&gt;I-&gt;A-&gt;I-&gt;A</b>
<b>CIFAR 10 -Quick</b>	<b>Object Recognition</b>	<b>C-&gt;P-&gt;A-&gt;N-&gt;C-&gt;A-&gt;P-&gt;N-&gt;C-&gt;A-&gt;P-&gt;I</b>
<b>CIFAR 10 -Full</b>	<b>Object Recognition</b>	<b>C-&gt;P-&gt;A-&gt;C-&gt;A-&gt;P-&gt;C-&gt;A-&gt;P-&gt;I-&gt;I</b>
<b>DJINN ASR</b>	<b>Speech to text Decoder</b>	<b>I-&gt;A-&gt;I-&gt;A-&gt;I-&gt;A-&gt;I-&gt;A-&gt;I-&gt;A-&gt;I-&gt;A-&gt;I</b>

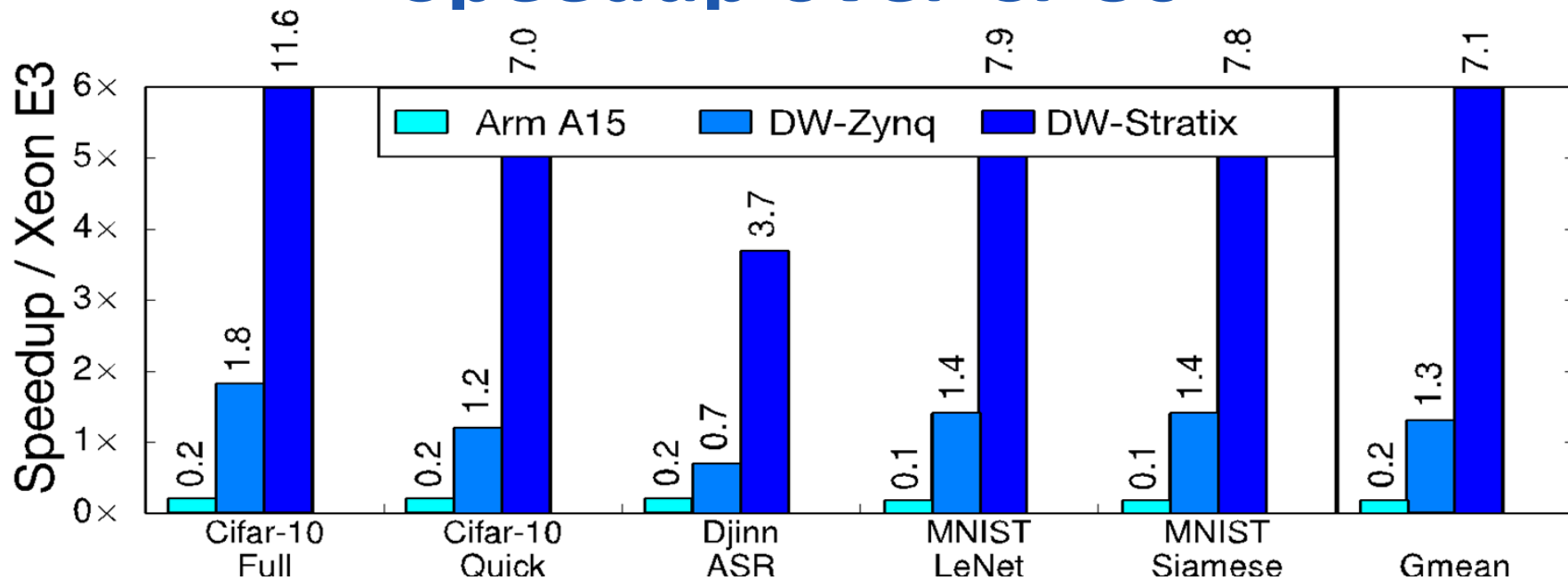


# Evaluated Platforms

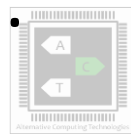
<b>FPGA</b>	<b>Altera Stratix V</b> TDP:25W \$6999	<b>Xilinx Zynq 7000 ZC702</b> TDP: 2W \$129	
<b>CPU</b>	<b>Intel Xeon E3-1276 V3</b> TDP: 84W \$339	<b>ARM Cortex 15</b> TDP: 5W \$191	
<b>GPU</b>	<b>Tegra K1 GPU</b> TDP: 10 W \$191	<b>GeForce GTX 650 Ti</b> TDP: 110 \$150	<b>Tesla K40</b> TDP: 235 W \$5499



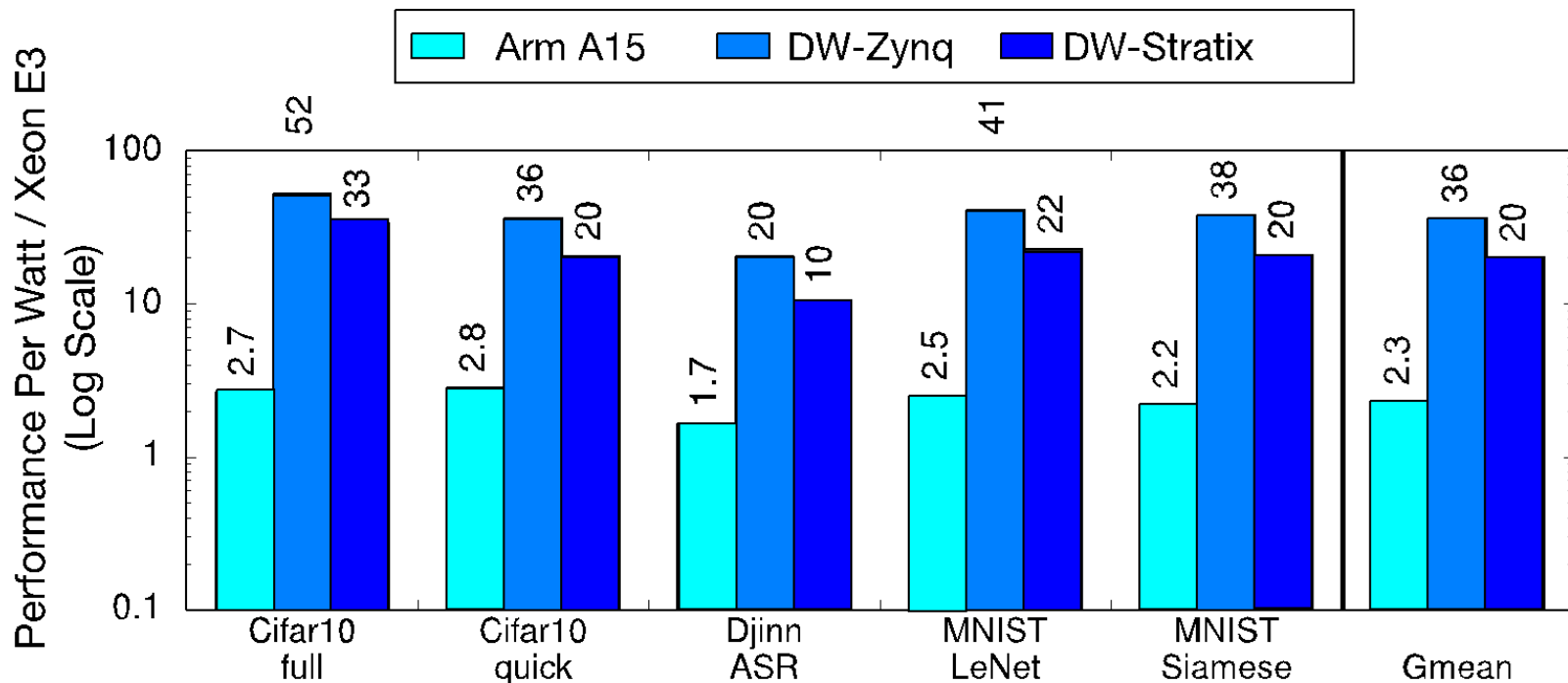
# Speedup over CPUs



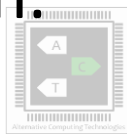
DnnWeaver achieves a speedup of up to **7.1x** Xeon.



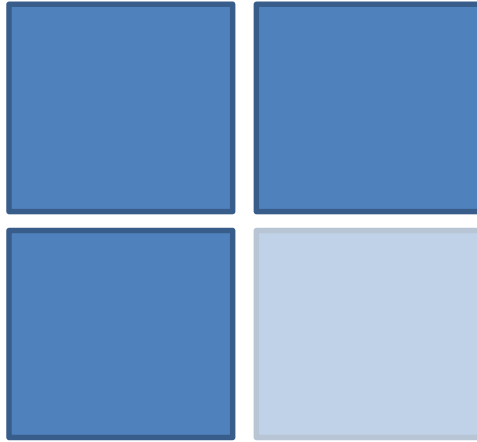
# Performance-per-watt over CPUs



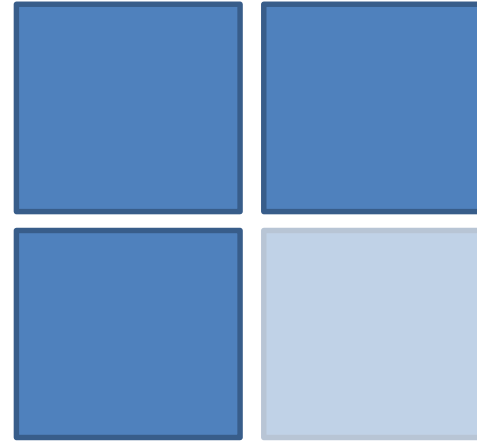
DnnWeaver is up to **36x** more power efficient than Xeon.



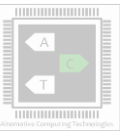
# Conclusion



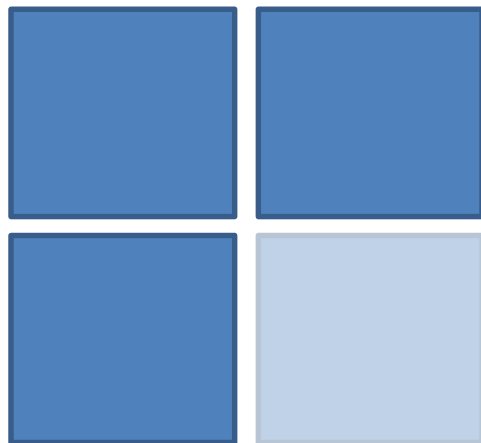
3 cores available



3 cores available

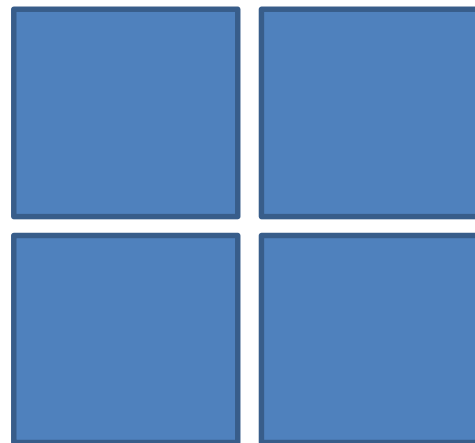


# Conclusion

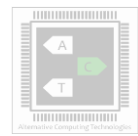


3 cores available

## Microfluidic Cooling



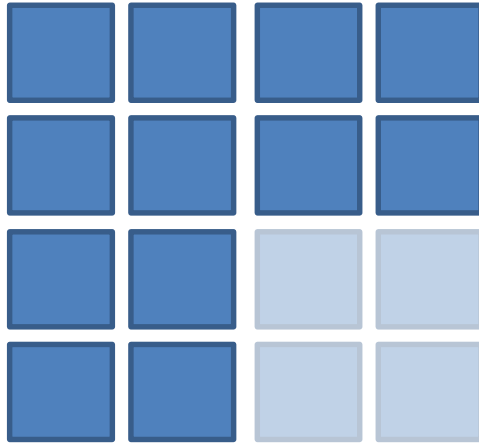
4 cores available





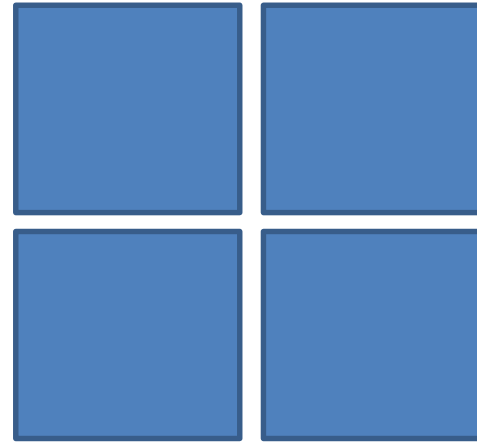
# Conclusion

## Approximate Computing

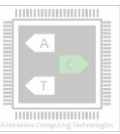


12 cores available

## Microfluidic Cooling

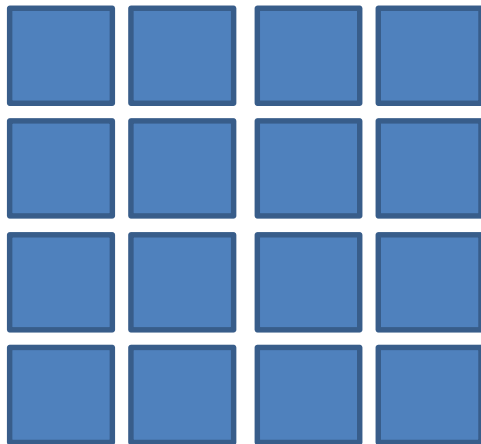


4 cores available

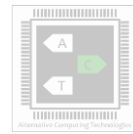


# Conclusion

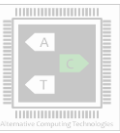
**Approximate Computing**  
+  
**Microfluidic Cooling**



16 cores available

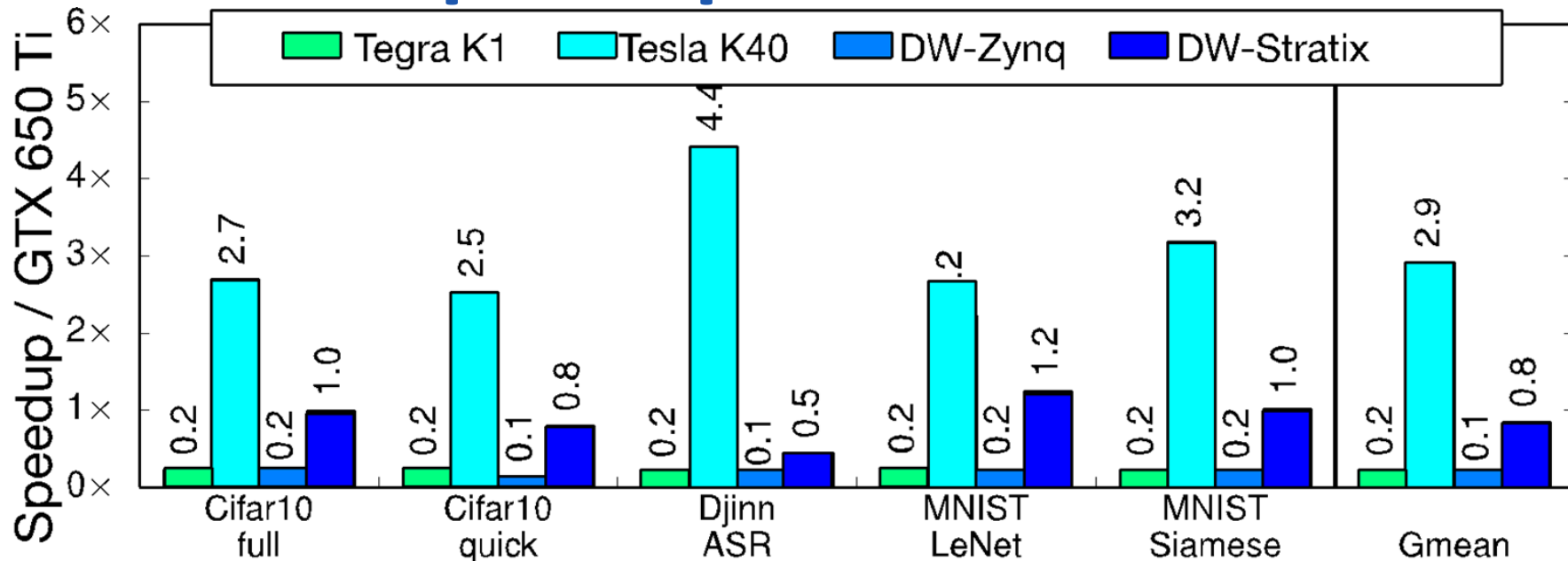


# Questions?

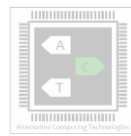


Backup

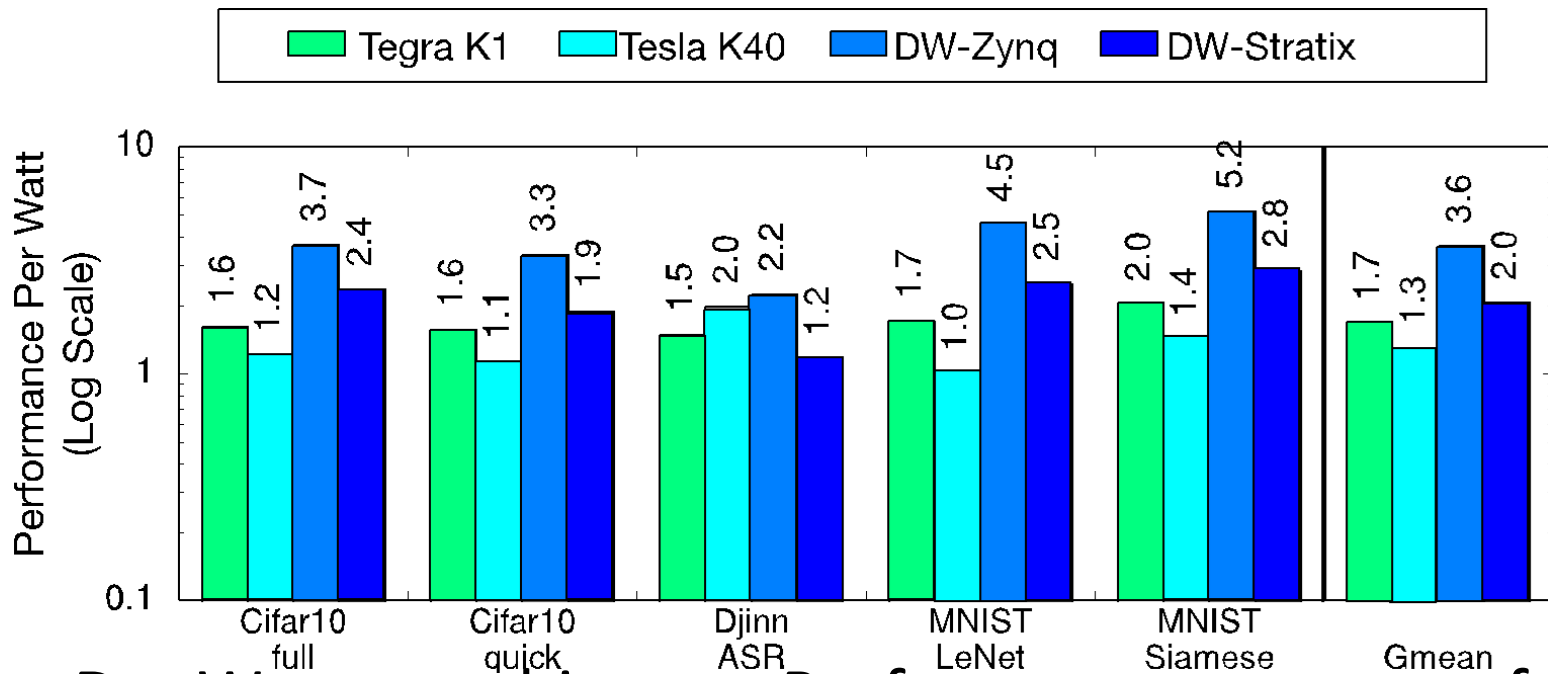
# Speedup over GPUs



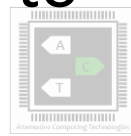
DnnWeaver provides an average of **0.8x** speedup over GTX 650Ti



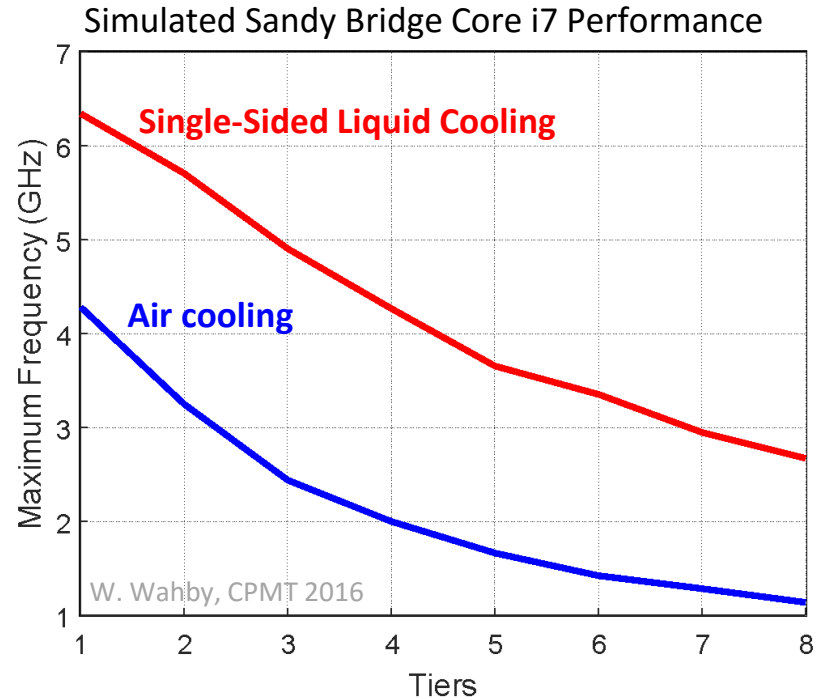
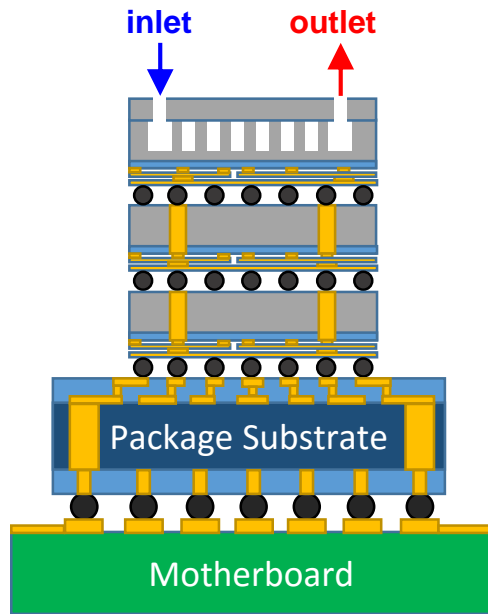
# Performance-per-watt over GPUs



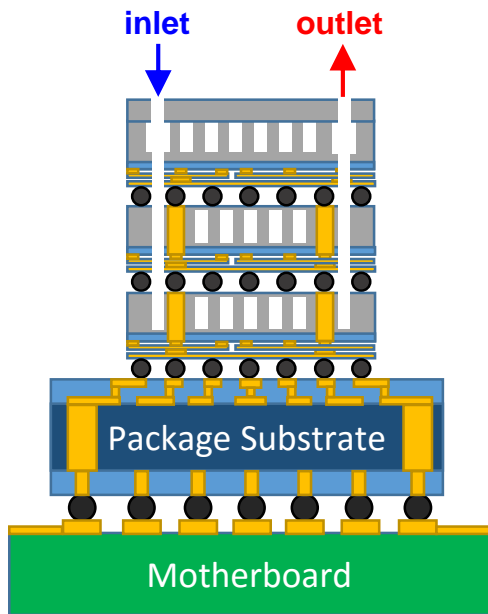
DnnWeaver achieves a Performance per watt of up to  
**3.6x** GTX 650 Ti



# 3D stacking complicates heat removal



# Tier-specific microfluidic cooling



Simulated Sandy Bridge Core i7 Performance

