

Arm **2018**

A Vision for Destruction of Post-Moore Disruption: How I Learned To Stop Worrying and Love Non-Von Neumann

Presented by: Jonathan Beard

UCAR SEA, Boulder, CO

3 April

Twitter: @jonathan_beard

Copyright © 2018 Arm Limited

Disclaimer

- This is not a product presentation.
- Any changes to architecture must go through the standard architecture review board process.
- Nothing within this presentation should be taken as a direction of future architecture or product plans for Arm Inc/Ltd.
- Please ask questions related to the presentation, this should be a fun presentation!!

About Me



Pale Peppered Moth



src: https://goo.gl/QJbLqG

19th Century London



src: https://goo.gl/BA72Gu

Bird Food



arm

Pale Peppered Moth



src: https://goo.gl/QJbLqG



Pale Peppered Moth



src: <u>https://goo.gl/1BBZVG</u>

Pale Peppered Moth



src: https://goo.gl/QJbLqG





src: https://goo.gl/1BBZVG



src: https://goo.gl/yna3tH

Pale Peppered Moth



src: https://goo.gl/QJbLqG



Selection on Computers



src: https://goo.gl/1dQoCZ

Speed \$\$ **Energy Consumption** Ease of Use

Software Availability















Source: Shekhar Borkar, Journal of Lightwave Technology, 2013

Selection on Computers



src: https://goo.gl/1dQoCZ

Speed SS Energy Consumption

Ease of Use Software Availability







Evolutionary Pressure



Evolutionary Pressure



Many Language Options

Java StormObj-CR Go SwiftOpenMPScala X10 FortranPerl C++Javascript Ruby PythonChapel Rust CPascal Spark MPI





Why language interfaces stabilize







Evolutionary Pressure





Accelerator



Hyper-optimized pressure

Accelerator



Orm

Coevolution



src: https://goo.gl/1dQoCZ



_MM_TRANSPOSE4_PS(*((__v4sf*)b[4]) , *((__v4sf*)b[5]) , *((__v4sf*)b[6]) , *((__v4sf*)b[7])); _MM_TRANSPOSE4_PS(*((__v4sf*)(b[4]+4)) , *((__v4sf*)(b[5]+4)) , *((__v4sf*)(b[6]+4)) , *((__v4sf*)(b[7]+4))); row_swap((__v4sf*)(b[0]+4),(__v4sf*)(b[4])); row_swap((__v4sf*)(b[1]+4),(__v4sf*)(b[5])); row_swap((__v4sf*)(b[2]+4),(__v4sf*)(b[6])); row_swap((__v4sf*)(b[3]+4),(__v4sf*)(b[7]));

orm

Coevolution

```
_MM_TRANSPOSE4_PS(*((__v4sf*)b[4]) , *((__v4sf*)b[5]) , *((__v4sf*)b[6]) , *((__v4sf*)b[7]));
_MM_TRANSPOSE4_PS(*((__v4sf*)(b[4]+4)) , *((__v4sf*)(b[5]+4)) , *((__v4sf*)(b[6]+4)) , *((__v4sf*)(b[7]+4)));
row_swap((__v4sf*)(b[0]+4),(__v4sf*)(b[4]));
row_swap((__v4sf*)(b[1]+4),(__v4sf*)(b[5]));
row_swap((__v4sf*)(b[2]+4),(__v4sf*)(b[6]));
row_swap((__v4sf*)(b[3]+4),(__v4sf*)(b[7]));
```



src: https://goo.gl/1dQoCZ

Coevolution



- At some point we have to break out of the cycle.
- Cost of not doing so is getting more of the exact same architecture.
- Perhaps worse, architectures will evolve without much application input.





Some disruption will minimize itself

Stabilizing selection (technical debt)

- Keeps rapid oscillations in hardware stack under control
- Convergent evolution driven by software to hardware will drive similar solutions from multiple vendors (also potentially a bad thing...)
- How do we minimize cost of evolutionary changes while maximizing ability to introduce variation?

Neuromorphic Accelerator Quantum Accelerator Compute In Memory (computational

SRAM)

FPGA

Programming likely won't change much

Pale Peppered Moth



src: <u>https://goo.gl/QJbLqG</u>





src: https://goo.gl/1BBZVG



src: https://goo.gl/yna3tH

Pale Peppered Moth



src: https://goo.gl/QJbLqG

How to minimize disruption

Pale Peppered Moth



src: https://goo.gl/1BBZVG

Pale Peppered Moth



src: https://goo.gl/QJbLqG





34 Confidential © Arm 2018

arm

Common core: multi-part chip





Minimize the cost of being eaten

- Enable faster integration of accelerators
- Enable easier integration of new memory technologies
- Enable code to run on multiple architectures more transparently
- Streamline virtual memory and security so accelerators are included natively
- Minimize data movement (tomorrow's talk)
- Most importantly: re-enable portability



Folding (a.k.a. Manual Virtual Memory)

- Virtual memory started b/c of lack of available memory vs. storage
- Programmers wrote code that manually folded/unfolded to storage
- Huge controversy existed b/c all code at the time was written this way – Automatic folding hardware was ~25% slower than manual folding



"In ceasing to expend energy (item 1) in a process whose main result is to make programs less fit to run on other machine configurations (item 2), or to run in company with other programs (item 3), or to run with temporarily reduced resources (item 4), we do more than reduce costs; we remove self-created obstacles which today are impeding the development of needed types of systems"

- D. Sayre, IBM Yorktown Research (1969)

1970 - The Circle of Tech - "Self Created Obstacles"



Constraining force: cost of extending virtual memory



Memory Interface

Software Interface —? ———— Hardware Interface





Example: current context swap behavior

- Newer NV tech enables access in 100-2000 cycle range
- Reduce context swap overhead (currently Linux overhead for swap is 2000—5000 cycles) to enable hiding of intermediate latencies.



Secondary context swap behavior





Conclusion

- > We will see some serious attempts at disruption over next 5-20 years
- > Many disruption will die quickly b/c of our technical debt
- Providing a better hardware/software abstraction can enable faster evolution while reducing destabilizing impact to software of technology disruption.
- Better interfaces are needed for:
 - Virtual memory
 - Context migration
 - Exception handling
- Slides will be available at: <u>http://www.jonathanbeard.io/media</u>
- Feel free to Tweet me @jonathan_beard
- > Thanks for listening!

Confidential © Arm 2018