# Amber SoC Tapeout Overview
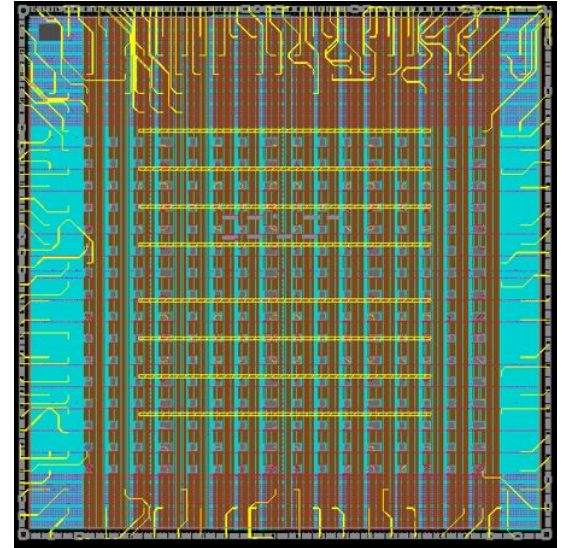
Stanford University -- AHA Monthly Meeting

December 3, 2020
Christopher Torng, Taeyoung Kong
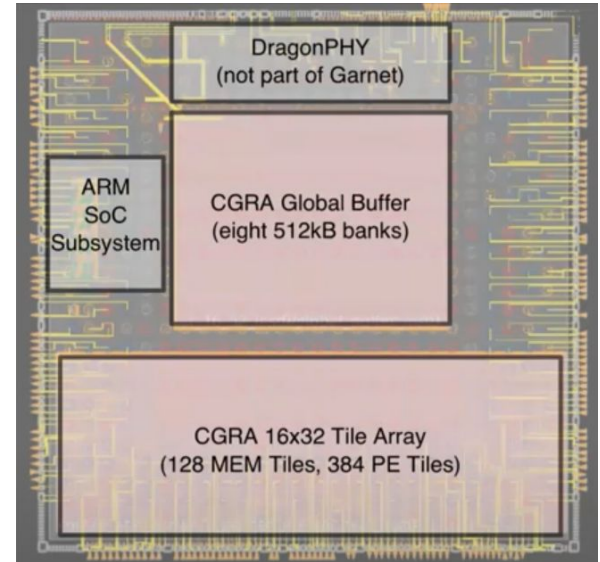
AHA
Agile Hardware Center

# AHA SoC Evolution -- Jade CGRA

- 16x16 array of PE and memory tiles

- Simple 16b fixed point PE tiles

- Memory tiles with line buffer access pattern only

- Image processing applications



AHA
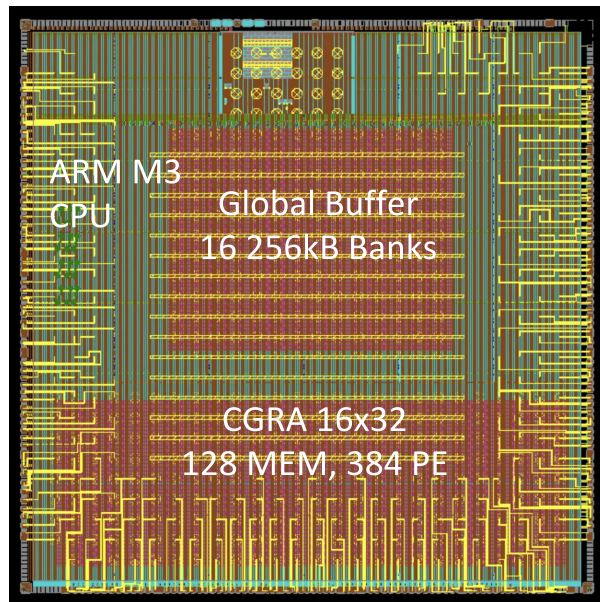Agile Hardware Center

# AHA SoC Evolution -- Garnet SoC

- Virtually taped out on June 2020

- Jade PE + **FP multiply and add (BFloat16)**

- Memory tiles + **general affine access pattern**

- Image processing + ML

# AHA SoC Evolution -- Amber SoC

- Taped out on December 1, 2020

- Garnet PE + **Local register file (ponds)**

- Garnet memory tiles + **Easier to program from the compiler**

- Verified fast reconfiguration, application partitioning, temporal and spatial multiplexing of applications

- Image processing and ML



ARM M3 CPU
Global Buffer
16 256kB Banks
CGRA 16x32
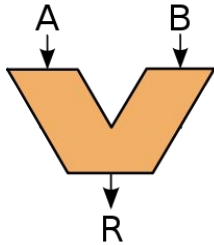128 MEM, 384 PE

AHA
Agile Hardware Center

# Amber SoC Architecture
## Memory System for Virtualization and Energy-Efficiency

# Tape-out Revisions and Memory System

| Tapeout Revision | Hierarchy | Features |
|---|---|---|
| Jade<br>(Gen-1 CGRA) | Dram -><br>Memory Tile | First version. No SoC |
| Garnet<br>(Gen-2 CGRA + SoC) | Dram -><br>Global Buffer -><br>Memory Tile | Energy-efficiency with 3-level memory system |
| Amber SoC<br>(Gen-3 CGRA + SoC) | Dram -><br>Global Buffer -><br>Memory Tile -><br>Pond (Compute-local RF) | **Virtualization**<br>Energy-efficiency with 4-level memory system |

AHA
Agile Hardware Center

# Introduction to CGRA Virtualization

- CGRA virtualization at different levels



Block-level



CGRA-level



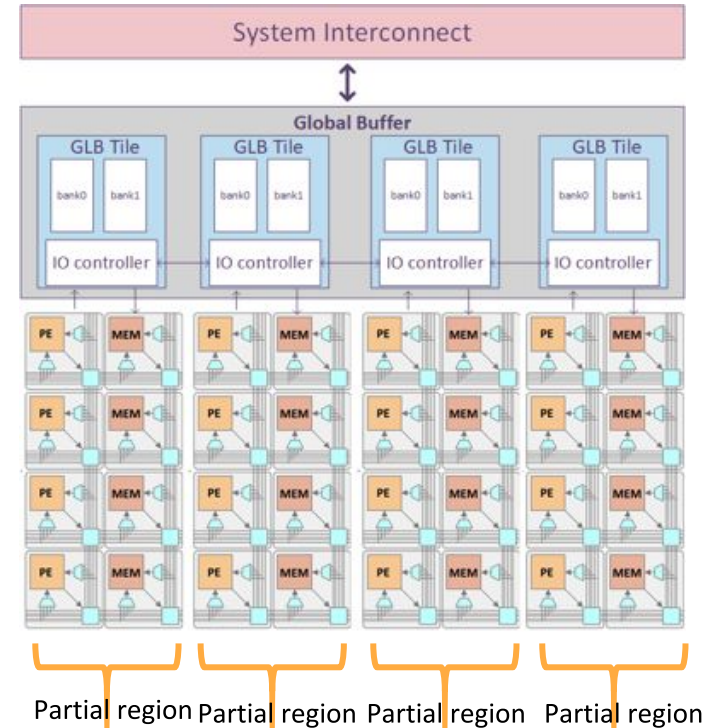Multi node level

- Definition / Advantages
  - Abstract CGRA hardware resources
  - Simplify the interface
  - Support flexible execution models

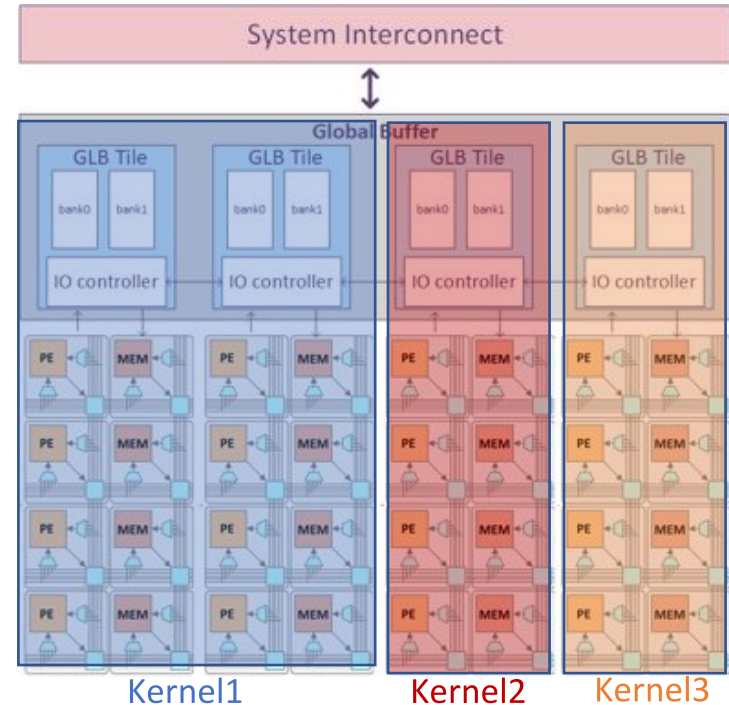=> Requires hardware and software support

# Amber: An Initial Prototype with Hardware Mechanisms for Virtualization

- Mechanism #1 : Run multiple kernels in parallel

- Mechanism #2: Relocatable bitstream to different partial regions

- Mechanism #3: Dynamic partial reconfiguration



Partial region   Partial region   Partial region   Partial region
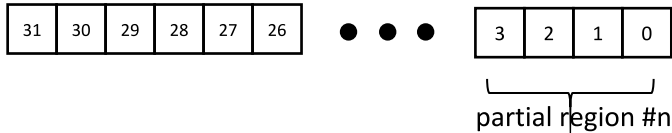
# Hardware Mechanism #1 – Multiple Kernels in Parallel

- Three kernels running independently and simultaneously
  - Virtualization at I/O, GLB, and Tile

- I/O allocation
  - Not limited to the location of the memory bank

- memory/computation allocation
  - MEM: Different # of global buffer
  - COMP: Different # of partial regions



AHA
Agile Hardware Center

# Hardware Mechanism #2 – Relocatable Bitstream to Different Partial Regions

● Bitstream address breakdown

| 31 | 30 | 29 | 28 | 27 | 26 | ● ● ● | 3 | 2 | 1 | 0 |
|----|----|----|----|----|----|-------|---|---|---|---|

partial region #n

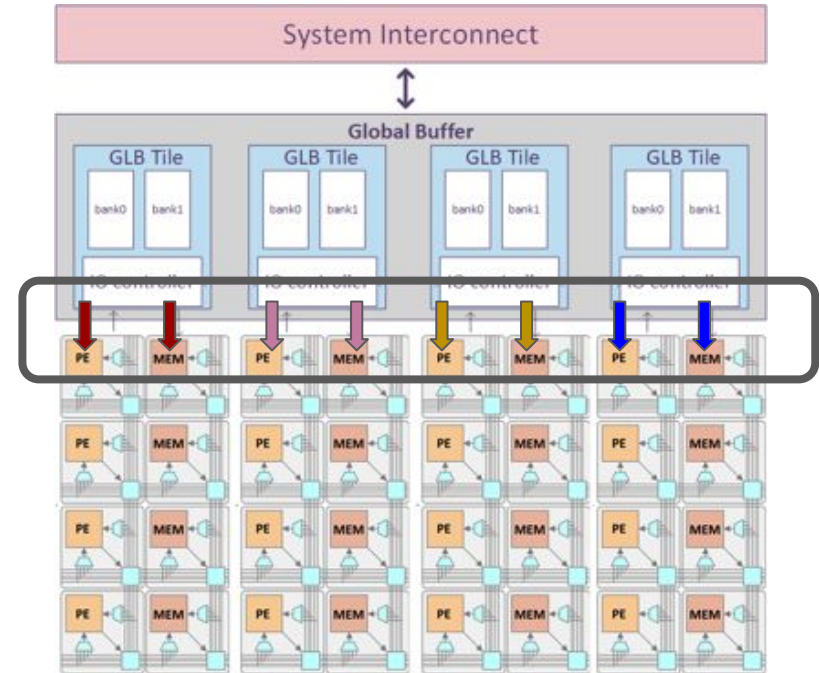- ○ LSB 8bits indicate the column #
- ○ Rest address bits and data bits remain unchanged



Relocatable

AHA
Agile Hardware Center

# Hardware Mechanism #3 – Dynamic Partial Reconfiguration (DPR)

- Use global buffer to double buffer bitstream
  - Fast
  - Parallel
  - Independent

- Configuration time (75,520 registers)
  - Jtag: 108.7ms
  - AXI4-Lite: 604us
  - DPR: 4.72us



AHA
Agile Hardware Center

# Tape-out Revisions and Memory System

| Tapeout Revision | Hierarchy | Features |
|---|---|---|
| Jade<br>(Gen-1 CGRA) | Dram -><br>Memory Tile | First version. No SoC |
| Garnet<br>(Gen-2 CGRA + SoC) | Dram -><br>Global Buffer -><br>Memory Tile | Energy-efficiency with 3-level memory system |
| Amber SoC<br>(Gen-3 CGRA + SoC) | Dram -><br>Global Buffer -><br>Memory Tile -><br>Pond (Compute-local RF) | Virtualization<br>**Energy-efficiency with 4-level memory system** |

AHA
Agile Hardware Center

# Hardware Support for Energy-Efficient Memory System

- More energy-efficient memory system
  - Reduce memory access
  - Memory hierarchy
  - Unified framework for each level of hierarchy



AHA
Agile Hardware Center

# Hardware Support for Energy-Efficient Memory System

- More energy-efficient memory system
  - Reduce memory access
  - Memory hierarchy
  - Unified framework for each level of hierarchy

- Mechanism 1: Leverage wide-fetch, single-port SRAM for better efficiency

# Hardware Support for Energy-Efficient Memory System

- More energy-efficient memory system
  - Reduce memory access
  - Memory hierarchy
  - Unified framework for each level of hierarchy

- Mechanism 1: Leverage wide-fetch, single-port SRAM for better efficiency

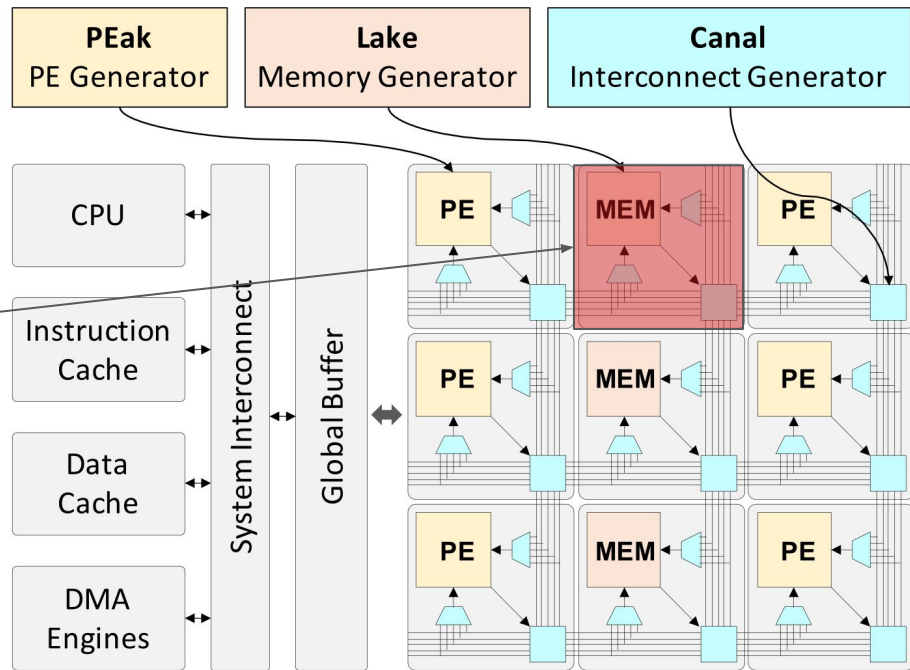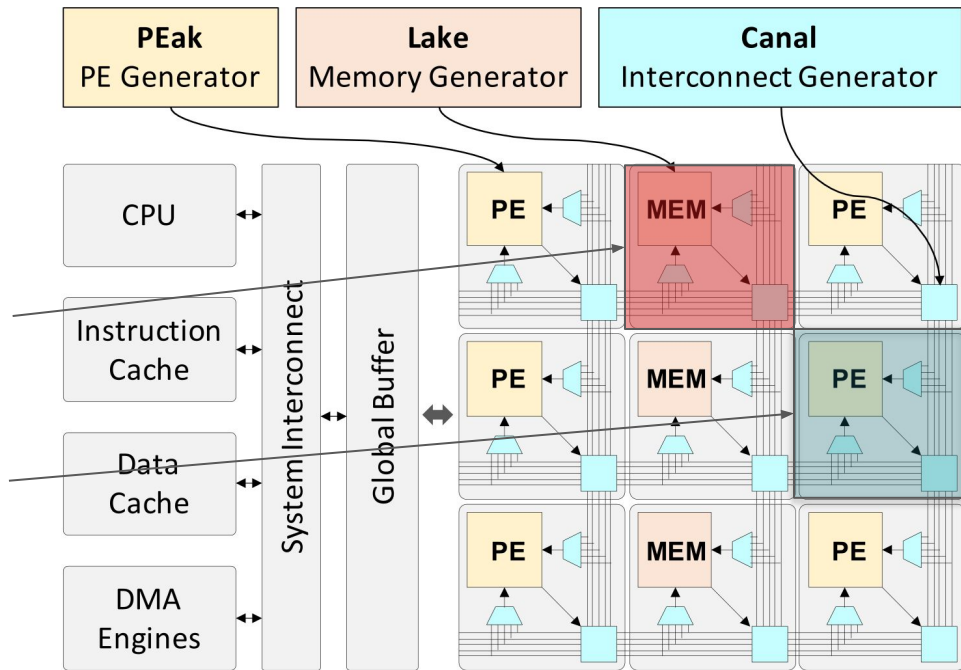- Mechanism 2: Place small register-based memories in the PE

# Hardware Support for Energy-Efficient Memory System

- More energy-efficient memory system
  - Reduce memory access
  - Memory hierarchy
  - Unified framework for each level of hierarchy

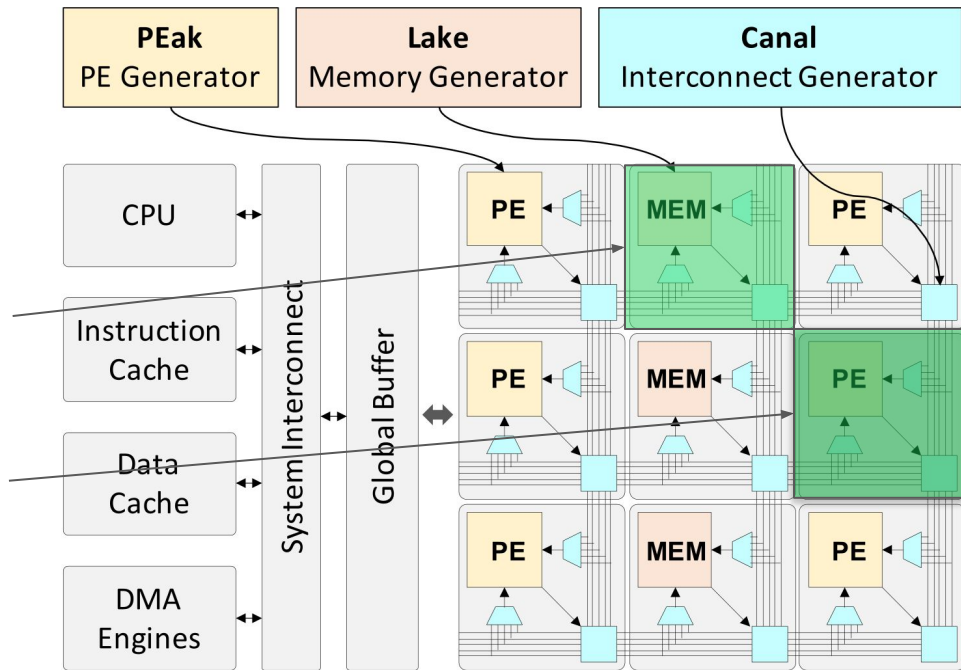- Mechanism 1: Leverage wide-fetch, single-port SRAM for better efficiency

- Mechanism 2: Place small register-based memories in the PE

- Mechanism 3: Unified memory generator framework for every level of hierarchy

# Hardware Mechanism #1 - Wide Fetch SRAM

- Mechanism 1: Leverage wide-fetch, single-port SRAM
  - Reduce memory read/write

- Wide-fetch SRAMs exhibit better energy per data accessed compared to single-wide SRAMs

- Wide-fetch SRAMs provide equivalent bandwidth with only a modest control overhead compared to Dual-Ported SRAMs.



AHA
Agile Hardware Center

# Hardware Mechanism #2 - PE-local Register File

- Mechanism 2: Register file in the PE
  - The lowest level memory hierarchy

- Keep data local to PE for accumulation and low-latency buffering

- Avoid SRAM accesses by exploiting locality

# Hardware Mechanism #3 - Unified Framework

So far, the discussion has concerned mainly architecture...what about control?



- Mechanism 3: Unified framework of memory generator
  - Use the same framework for every memory level in the hierarchy
  - Different SRAM organization and addressing complexity

- Compiler analyzes looping structures of the application with given resource constraints
  - Generate configuration registers

AHA
Agile Hardware Center

# Summary

- **Amber** is a prototype for hardware mechanisms for CGRA **virtualization**
  - Multi-kernels in parallel
  - Relocatable bitstream
  - Dynamic partial reconfiguration

- **Amber** has an **energy-efficient** multi-level memory system
  - Wide-fetch, single-port SRAM in MEM tile
  - Register-file in PE tile
  - Unified memory generator framework for each level of memory

AHA
Agile Hardware Center

# Tapeout Methodology

1. Building a chip in the cloud

2. Reflecting on a first attempt at agile physical development



ARM M3 CPU

Global Buffer
16 256kB Banks

CGRA 16x32
128 MEM, 384 PE

AHA
Agile Hardware Center

# We built Amber in TSMC Virtual Design Environment (VDE)

- **Vendor**: Prevent accidental leakage of TSMC's technology information
- **Customer**: Cloud EDA

# We built Amber in TSMC Virtual Design Environment (VDE)

Logging in from a browser

Access to a terminal

Display available for GUI needs



AHA
Agile Hardware Center

# VDE was a success in many respects

From TSMC and Cadence's perspective

- TSMC protected their sensitive information from accidental leakage
- Cadence has a new platform to offer as a service

From Stanford's perspective

- We **successfully built a chip** in the cloud
- We had access to all the tools, IP, and compute we needed
- Unique for academics -- We had access to live EDA support (Cadence)

**AHA**
Agile Hardware Center

# VDE introduced new pain points

| Challenge | Why it can hurt |
|---|---|
| Not allowed to extract design and source files from the closed environment (or difficult to) | Losing project history hurts, especially for open-source projects. Difficult to share code with non-VDE projects. |
| Cloud platform sometimes high latency | Interactive jobs can be very slow |
| Each VDE instance is a "new system" | Need to spend time getting tools to work again |
| Single-vendor toolchain (Cadence) | Less flexibility + Team may not have experts |
| Having to juggle cloud resources vs costs (e.g., making sure big jobs run on big servers) | Overpay for cloud resources or suffer the low performance |

AHA
Agile Hardware Center

# Our take on the VDE

- Modularity with **mflowgen helped us quickly adapt** between a single-vendor flow (Cadence) and a mixed-vendor flow (Synopsys, Cadence, Mentor)



| adk | design_merged.gds |
| cadence-pegasus-drc |
| drc.results | drc.summary |

swap

| adk | design_merged.gds |
| mentor-calibre-drc |
| drc.results | drc.summary |

| adk | design.pt.sdc | design.spef.gz | design.vcs.v |
| cadence-genus-genlib |
| design.lib |

swap

| adk | design.pt.sdc | design.spef.gz | design.vcs.v |
| synopsys-ptpx-genlibdb |
| design.db | design.lib |

- We maintained version control outside the VDE and **repeatedly uploaded our changes** (only lightweight changes)

- Other challenges decreased productivity (robustness, latency, tools, maintenance), but the **chip still went out successfully**
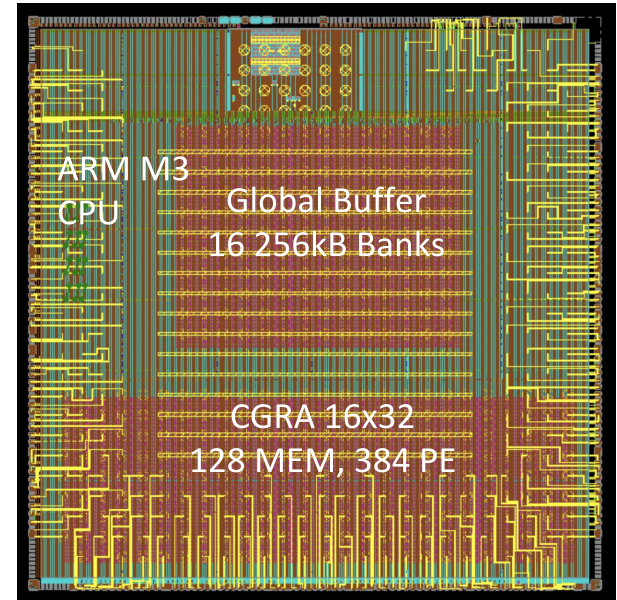
# Tapeout Methodology

1.  Building a chip in the cloud

    -   VDE was a success with room for improvement

2.  Reflecting on a first attempt at agile physical development

    -   Producing a golden tape every week



ARM M3 CPU
Global Buffer
16 256kB Banks
CGRA 16x32
128 MEM, 384 PE

AHA
Agile Hardware Center

# Weekly golden releases instead of "one big release"

Push away from "**one big release**" and towards software-like **frequent small ones**

**Big initial effort**

**Big effort**

**Big effort**

Move to VDE

Ponds

Memory fix

Timing fixes

**Big initial effort**

small

small

small

small

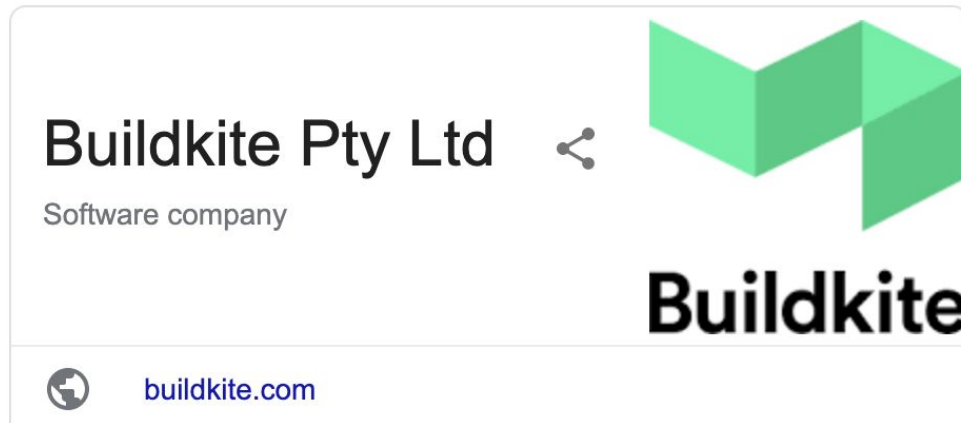# Weekly golden releases instead of "one big release"

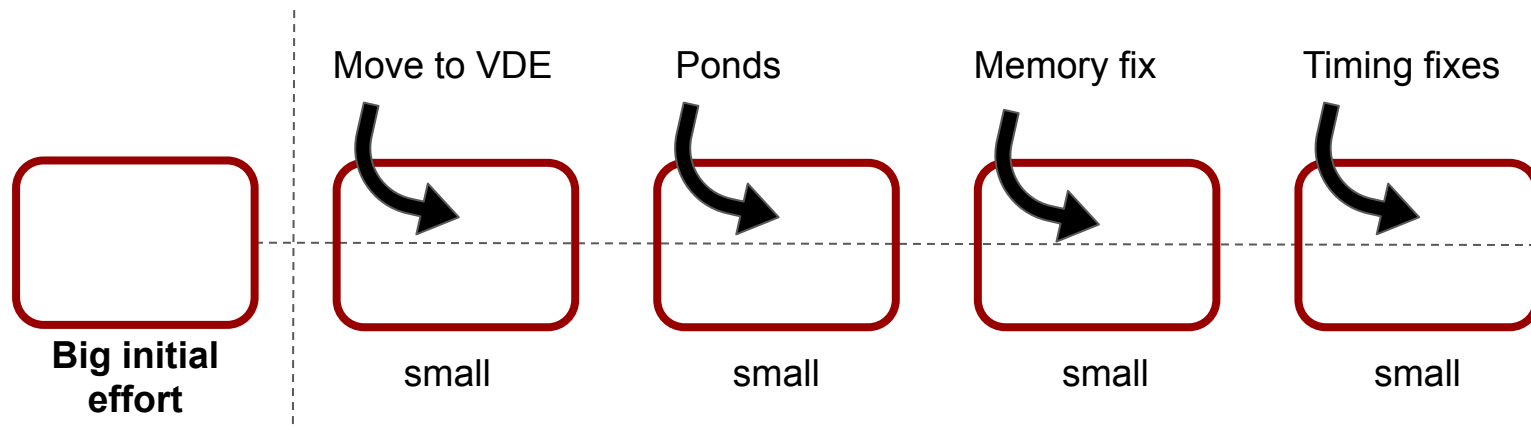- We were successful to a large degree...



```
% ls gold.archives

gold.112
gold.148
gold.158
gold.174
gold.1987
gold.30
gold.40
gold.84
gold.93
gold.README
```

A history of golden builds throughout design development

- Required **culture change** and also surfaced **fundamental challenges**

AHA
Agile Hardware Center

# Agile culture shift -- Biweekly sprints for team focus



Move to VDE          Ponds          Memory fix          Timing fixes

**Big initial
effort**          small          small          small          small

Biweekly sprints helped us convert many unordered parallel efforts into a **progression of team efforts** that "commit" one by one into the chip

AHA
Agile Hardware Center

# Agile culture shift -- Adapting Scrums for COVID



Geekbot

**Asynchronous scrums** completely replaced synchronous ones to great success
(e.g., staying in sync, text records, less meeting fatigue)

# Challenge #1 -- Some design approaches are friendlier than others



**Logical View**

**Physical View (complex constraints)**
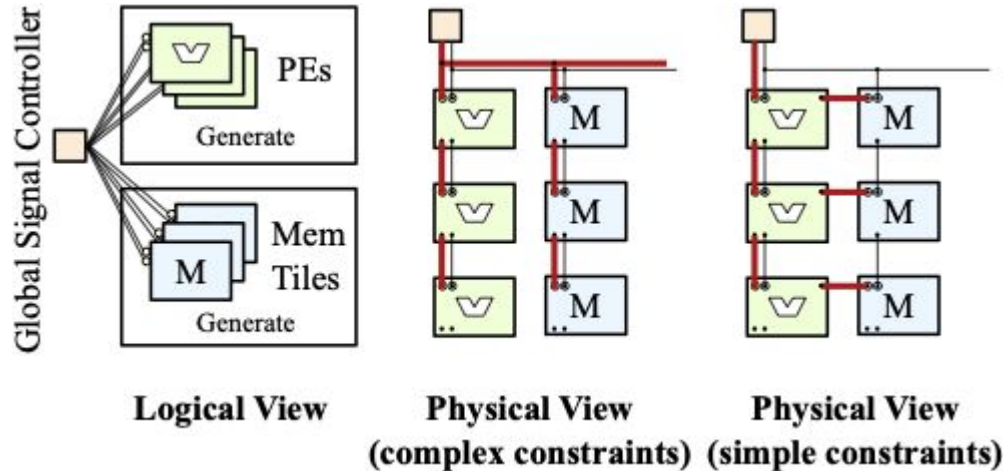
**Physical View (simple constraints)**

The (left) physical view requires **manual attention to constraints**.

The (right) physical view requires **negligible human attention**

While difficult to spot, certain design decisions which have no impact on logical design can cost weeks of physical iteration.
These need to be eliminated.

AHA
Agile Hardware Center

# Challenge #2 -- No single source of truth between physical/logical



**Logical View**

**Physical View** (complex constraints)

**Physical View** (simple constraints)

Significant design churn when…

- The **logical design updates, but physical collateral is stale** (e.g., constraints)

- The logical design updates **non-functionally** but generates a different netlist which breaks physical scripts

Need a mechanism to better tie together related logical and physical design files

**AHA** Agile Hardware Center

# Can we address these two challenges?

- Challenge #1 -- Some design approaches are friendlier than others

- Challenge #2 -- No single source of truth between physical/logical

AHA
Agile Hardware Center

# Summary of Amber SoC tapeout

- Taped out on December 1, 2020

- Garnet PE + **Local register file (ponds)**

- Garnet memory tiles + **Easier to program from the compiler**

- Verified fast reconfiguration, application partitioning, temporal and spatial multiplexing of applications

- Lessons learned in the **cloud**

- Lessons learned for **agile design**



TSMC 16nm
Cadence VDE
5mm x 5mm

ARM M3 CPU

Global Buffer
16 256kB Banks

CGRA 16x32
128 MEM, 384 PE

AHA
Agile Hardware Center