



香港中文大學
The Chinese University of Hong Kong

EDA for AI Chip Designs

Bei Yu
Chinese University of Hong Kong
byu@cse.cuhk.edu.hk

July 30, 2023



Autonomous drive

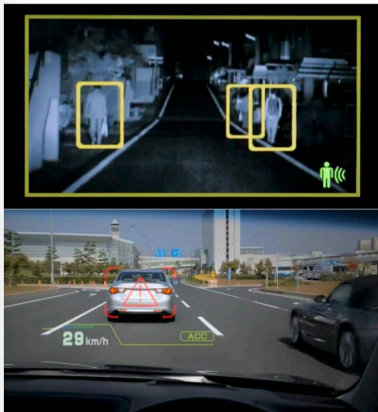
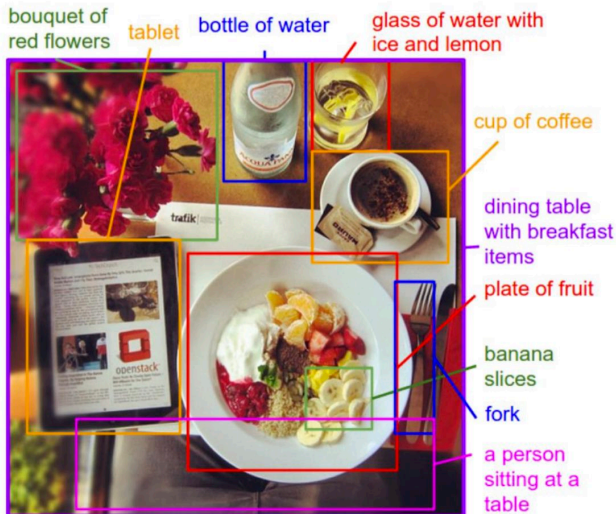
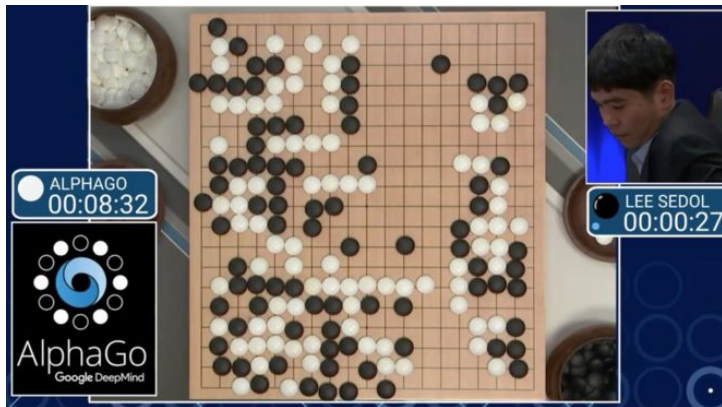


Image recognition

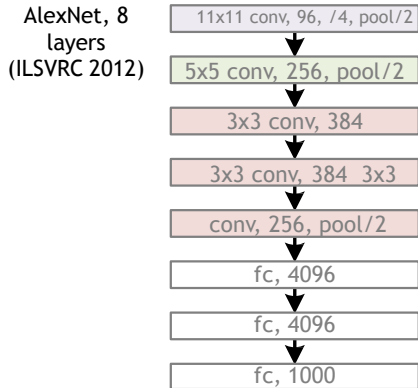




- A fast learning algorithm for deep belief nets. [Hinton et.al 1996]
- Data + Computing + Industry Competition
- NVidia's GPU, Google Brain (16,000 CPUs)
- **Speech**: Microsoft [2010], Google [2011], IBM
- **Image**: AlexNet, 8 layers [Krizhevsky et.al 2012] (26.2% -> 15.3%)



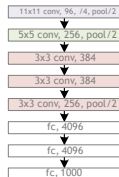
Revolution of Depth



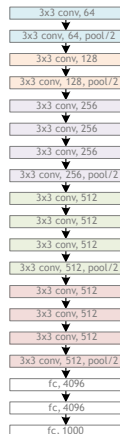
Slide Credit: He et al. (MSRA)

Revolution of Depth

AlexNet, 8 layers
(ILSVRC 2012)



VGG, 19 layers
(ILSVRC 2014)



GoogleNet, 22 layers
(ILSVRC 2014)



Slide Credit: He et al. (MSRA)

Revolution of Depth

AlexNet, 8
layers
(ILSVRC 2012)



VGG, 19
layers
(ILSVRC
2014)



ResNet, 152
layers
(ILSVRC 2015)



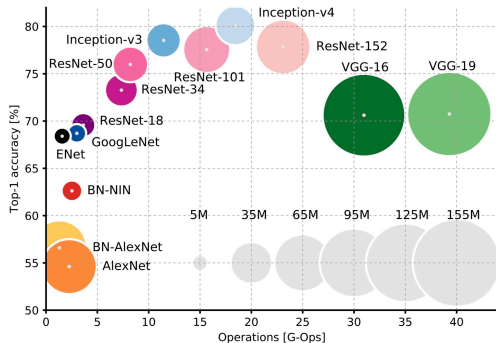
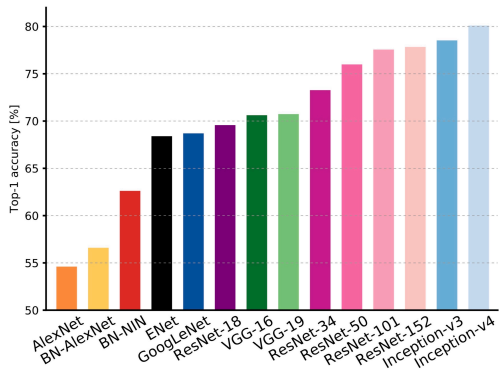
Slide Credit: He et al. (MSRA)



- AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB
- Network in Network (Lin, Chen, and Yan 2013) 29MB
- VGG (Simonyan and Zisserman 2015) 549MB
- GoogleNet (Szegedy, Liu, et al. 2015) 51MB
- ResNet (He et al. 2016) 215MB
- Inception-ResNet (Szegedy, Vanhoucke, et al. 2016)
- DenseNet (Huang et al. 2017)
- Xception (Chollet 2017)
- MobileNetV2 (Sandler et al. 2018)
- ShuffleNet (Zhang et al. 2018)

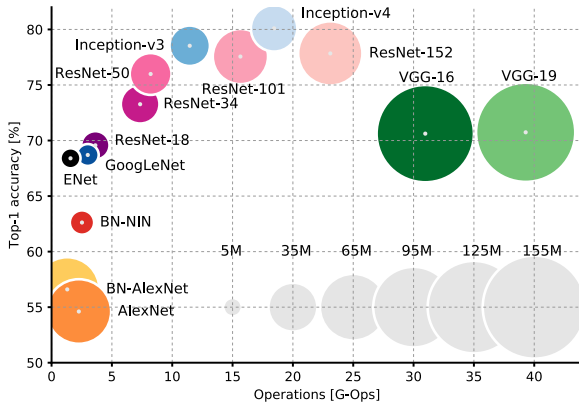


- AlexNet (Krizhevsky, Sutskever, and E. Hinton 2012) 233MB
- Network in Network (Lin, Chen, and Yan 2013) 29MB
- VGG (Simonyan and Zisserman 2015) 549MB
- GoogleNet (Szegedy, Liu, et al. 2015) 51MB
- ResNet (He et al. 2016) 215MB
- Inception-ResNet (Szegedy, Vanhoucke, et al. 2016) 23MB
- DenseNet (Huang et al. 2017) 80MB
- Xception (Chollet 2017) 22MB
- MobileNetV2 (Sandler et al. 2018) 14MB
- ShuffleNet (Zhang et al. 2018) 22MB



1

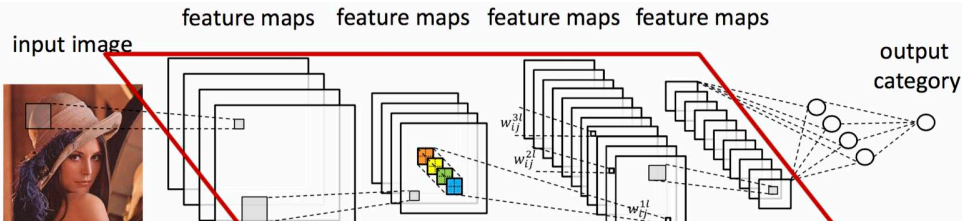
¹Alfredo Canziani, Adam Paszke, and Eugenio Culurciello (2017). “An analysis of deep neural network models for practical applications”. In: *arXiv preprint*.



Why AlexNet is large in size, but small in operations?

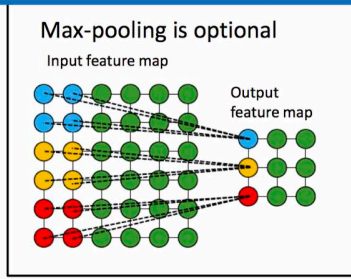
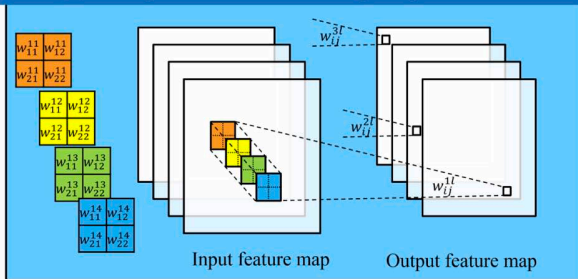
- Special FC layers
- Special Conv layers
- More channels
- Some redundant operators

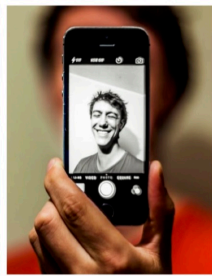
Convolutional Neural Network (CNN)



Convolutional layers account for over 90% computation

- [1] A. Krizhevsky, etc. Imagenet classification with deep convolutional neural networks. NIPS 2012.
- [2] J. Cong and B. Xiao. Minimizing computation in convolutional neural networks. ICANN 2014



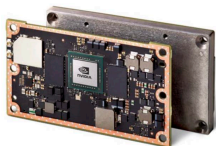


Embedded CV

Flexibility vs. Efficiency



CPU
(Raspberry Pi3)



GPU
(Jetson TX2)

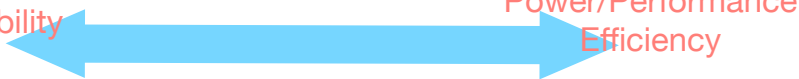


FPGA
(UltraZed)



ASIC
(Movidius)

Flexibility



Power/Performance
Efficiency

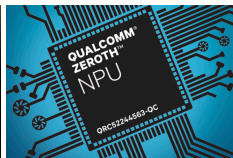
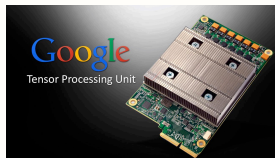


Convolution layer is one of the most expensive layers

- Computation pattern
- Emerging challenges

More and more end-point devices with limited memory

- Cameras
- Smartphone
- Autonomous driving

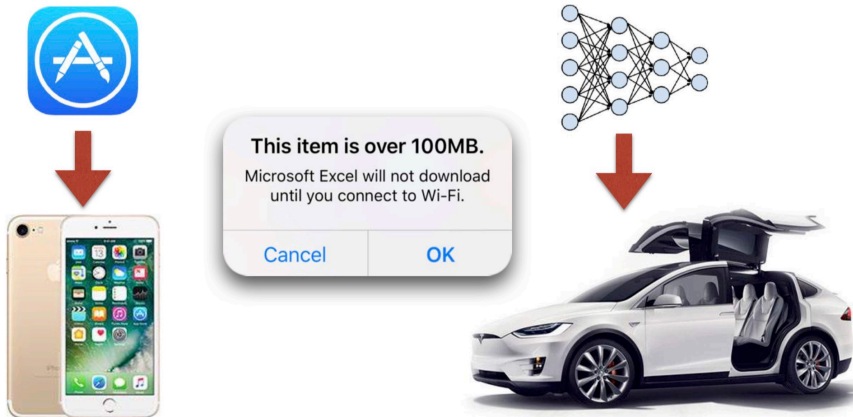


XILINX



An Intel Company

Hard to distribute large models through over-the-air update



2



AlphaGo: 1920 CPUs and 280 GPUs,
\$3000 electric bill per game



on mobile: **drains battery**
on data-center: **increases TCO**



3

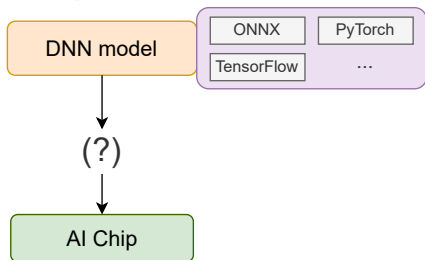
Application Category

Both	Datacenter	Edge
Intel, Nvidia, IBM, Xilinx, HiSilicon, Google, Baidu, Alibaba Group, Cambricon, DeePhi, Bitmain, Wave Computing	AMD, Microsoft, Apple, Tencent Cloud, Aliyun, Baidu Cloud, HUAWEI Cloud, Fujitsu, Nokia, Facebook, HPE, Thinkforce, Cerebras, Graphcore, Groq, SambaNova Systems, Adapteva, PEZY	Qualcomm, Samsung, STMicroelectronics, NXP, MediaTek, Rockchip, Amazon_AWS, ARM, Synopsys, Imagination, CEVA, Cadence, VeriSilicon, Videantis, Horizon Robotics, Chipintelli, Unisound, AISpeech, Rokid, KnuEdge, Tenstorrent, ThinCl, Koniku, Knowm, Mythic, Kalray, BrainChip, Almotive, DeepScale, Leepmind, Krtkl, NovuMind, REM, TERADEEP, DEEP VISION, KAIST DNP, Kneron, Esperanto Technologies, Gyrfalcon Technology, GreenWaves Technology, Lightelligence, Lightmatter, ThinkSilicon, Innogrit, Kortiq, Hailo, Tachyum

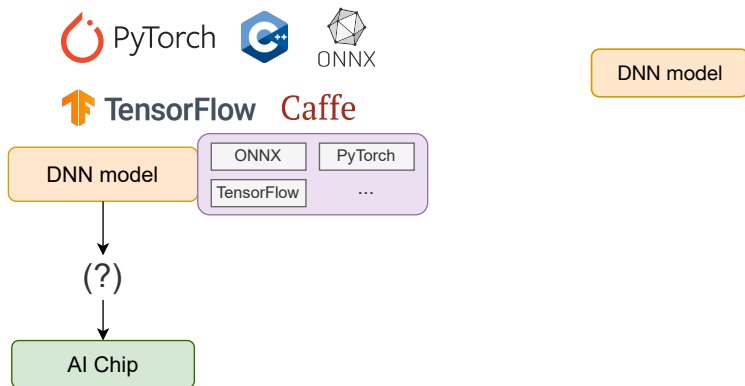
Source: <https://basicmi.github.io/Deep-Learning-Processor-List/>

DNN Deployment Flow

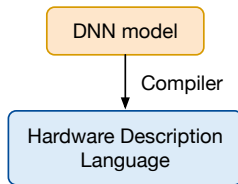
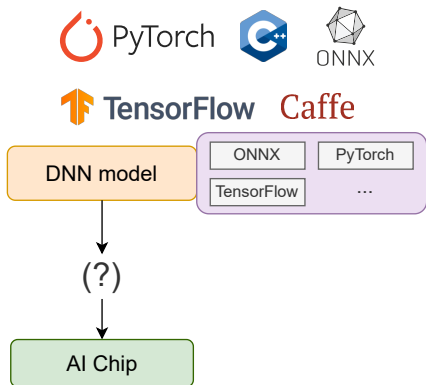
Deployment Flow: A Naive Approach



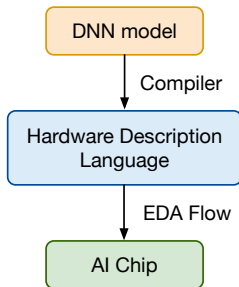
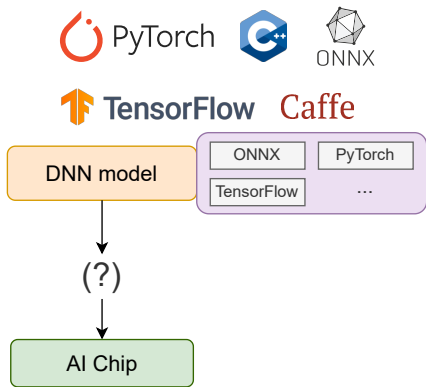
Deployment Flow: A Naive Approach



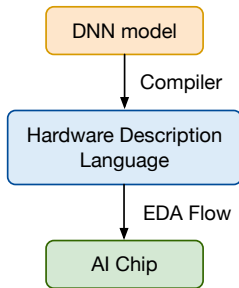
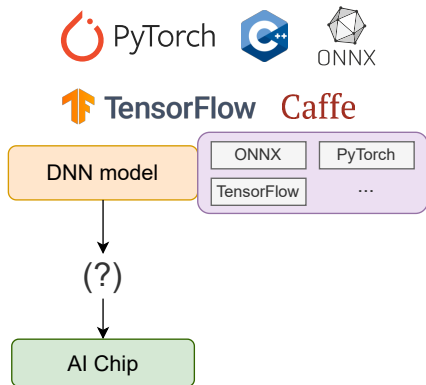
Deployment Flow: A Naive Approach



Deployment Flow: A Naive Approach



Deployment Flow: A Naive Approach

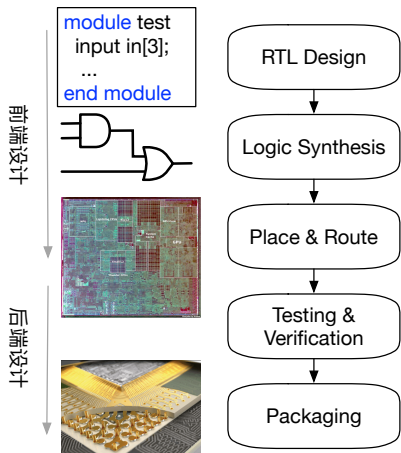
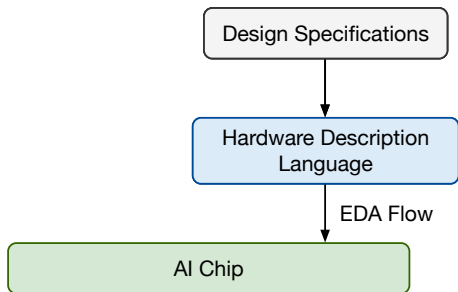


Question:

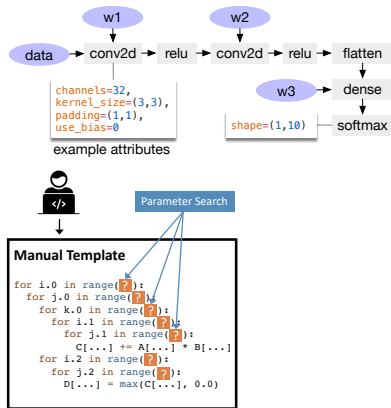
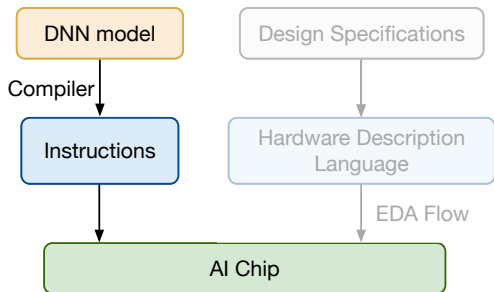
Why NOT using such deployment approach?



Deployment Flow: AI Chip Generation



Deployment Flow: DNN Model Compilation





LLVM

- Pro: Target-independent representation for optimization
- Con: Operating low-level operations is tricky

MLIR

- A tool for multi-level IR design (MLIR dialects)
- Enable different levels of abstraction

⁴<https://mlir.llvm.org/>



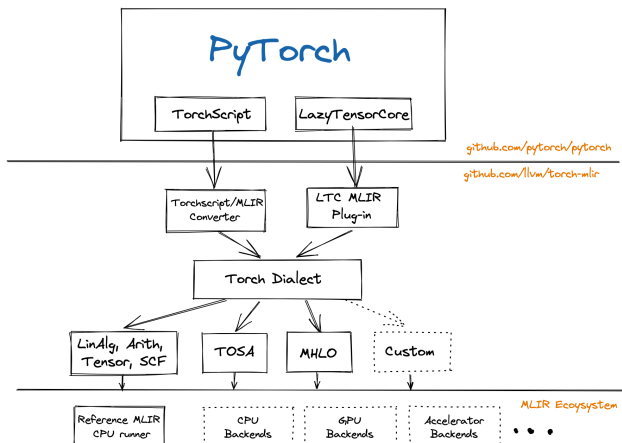
Torch-MLIR

The Torch-MLIR project aims to provide first class compiler support from the PyTorch ecosystem to the MLIR ecosystem.

Other MLIR Based DL Compiler:

- 1 OpenXLA
- 2 StableHLO
- 3 Triton
- 4 OneFlow

Torch-MLIR Architecture



The overall architecture of Torch-MLIR

⁵<https://github.com/llvm/torch-mlir>



TPU-MLIR⁶

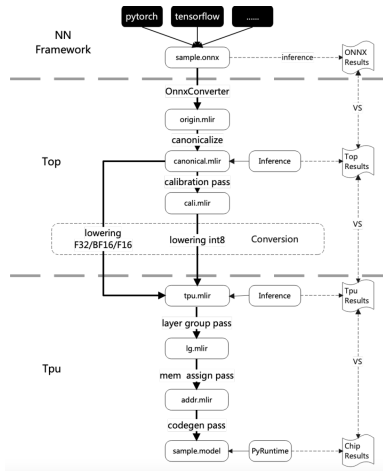
TPU-MLIR is an open-source machine-learning compiler based on MLIR for TPU.

Top Dialect:

- graph optimization
- quantization and inference
- ...

TPU Dialect:

- weight reordering
- operator slicing
- address assignment
- ...

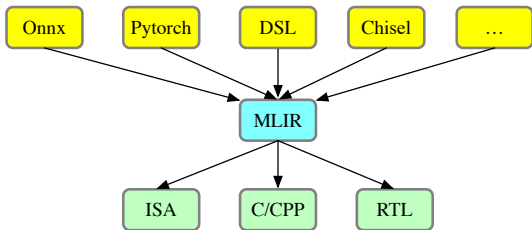


⁶Pengchao Hu et al. (2022). "TPU-MLIR: A Compiler For TPU Using MLIR". In: *arXiv preprint arXiv:2210.15016*.

Why MLIR?

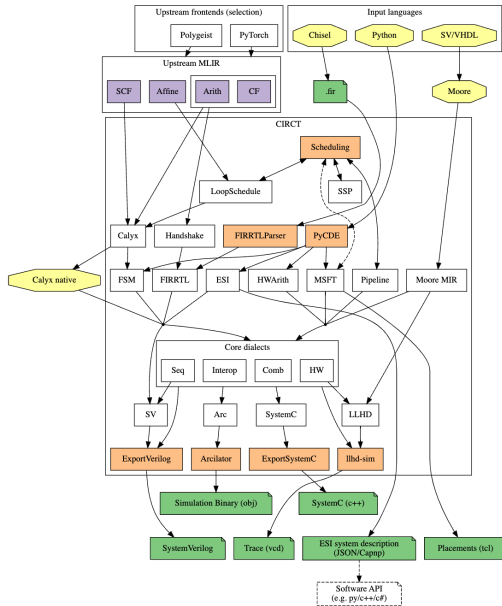
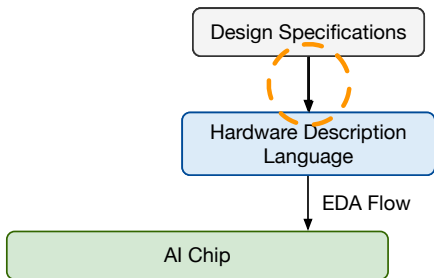
- Stable LLVM community
- Easy to implement
- Efficient optimization in multi-level abstract

★: You can simply combine other MLIR projects by dialect conversion.





- Circuit IR Compilers and Tools (CIRCT)
- apply MLIR and the LLVM development methodology
- to the domain of hardware design tools.





Calyx⁷

Calyx is a compiler infrastructure for languages that target hardware accelerators.

Firrtl⁸

Firrtl is an intermediate representation (IR) for digital circuits designed as a platform for writing circuit-level transformations.

Calyx and Firrtl has been integrated with the LLVM CIRCT infrastructure and is available as a dialect within it.



⁷Rachit Nigam et al. (2021). “A compiler infrastructure for accelerator generators”. In: *Proceedings of the 26th ACM International Conference on Architectural Support for Programming Languages and Operating Systems*, pp. 804–817.

⁸<https://github.com/chipsalliance/firrtl>

⁹<https://circt.llvm.org/>



Update on Chisel¹⁰

Chisel v3.6.0:

- The primary change in Chisel v3.6.0 is the transition from the Scala FIRRTL Compiler to the new MLIR FIRRTL Compiler(CIRCT/FIRTOOL).
- "Meaningless intermediate variable" problem is solved.

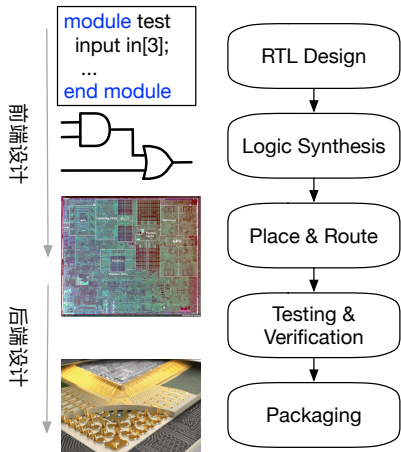
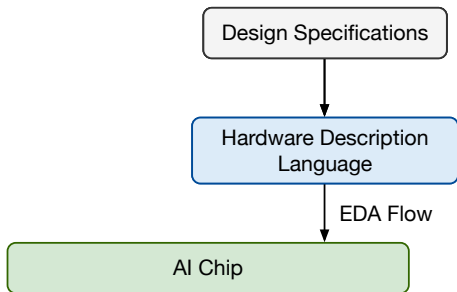
Chisel v5.0.0:

- Scala FIRRTL Compiler (SFC) is deprecated.



¹⁰Jonathan Bachrach et al. (2012). "Chisel: constructing hardware in a scala embedded language". In: *Proceedings of the 49th Annual Design Automation Conference*, pp. 1216–1225.

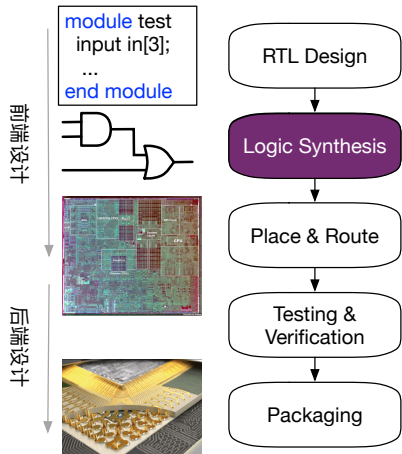
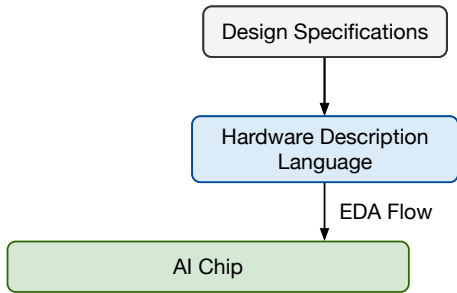
Today We Focus on EAD Part



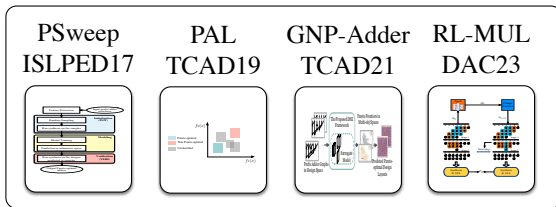


- 2 Arithmetic Unit Synthesis
- 3 Datapath Driven Placement
- 4 Wafer-Scale Floorplan

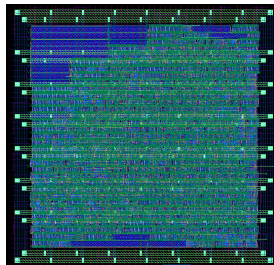
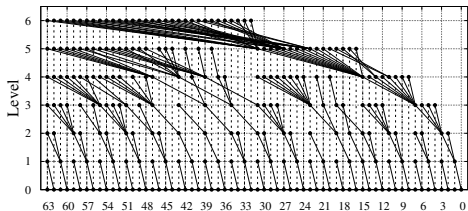
Arithmetic Unit Synthesis



Arithmetic Unit Synthesis



Logic synthesis v.s. Physical synthesis

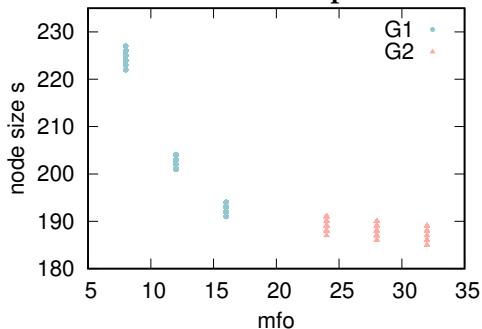


Constraints mapping between two synthesis stages is difficult.

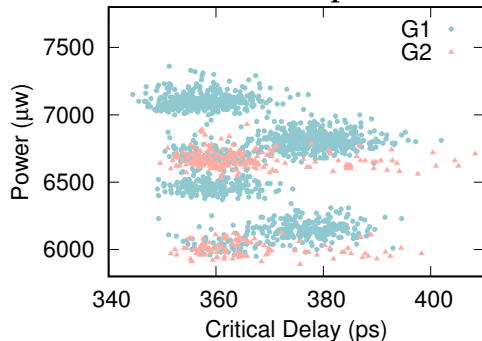
Question:

Why We Need to Optimize Arithmetic Circuits (Adder & Multiplier)?

Front-End Team Perspective:



Back-End Team Perspective:

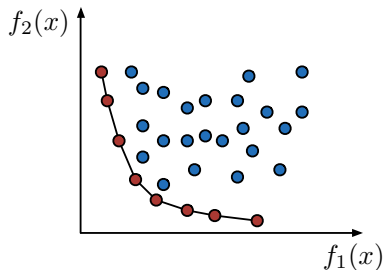


- Run design tools with all solutions is time-consuming. => **Enumeration is not feasible!**
- For 3K solutions, running time is $3000 \times 5 = 15\text{K}$ mins.
- What we care: **Pareto Frontier Curve**

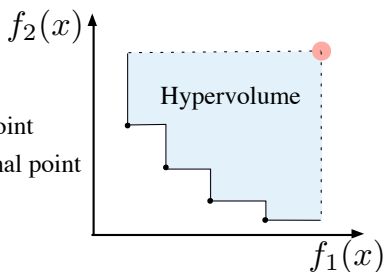


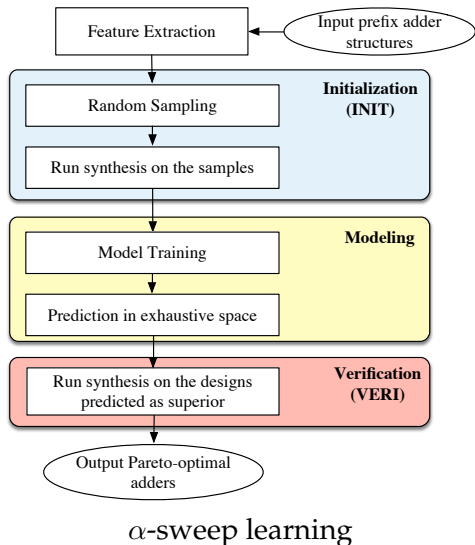
Pareto Frontier

- All the points are not dominated by any other point.
- Evaluation: Hyper-volume.
 - Size of the region bounded by the Pareto frontier and reference point.
 - Each dimension of the reference point is the maximum value on that dimension.



- Reference point
- Pareto-optimal point





- General flow;
- Use a joint output Power-Delay function (PD) as the regression output rather than using any single output;
- Select different features in different applications.

⁴S. Roy, Y. Ma, J. Miao and B. Yu, "A learning bridge from architectural synthesis to physical design for exploring power efficient high-performance adders", ISLPED'17.



Directions

- Can we improve the quality of the Pareto-Frontier?
- We use 3K samples to cover the solution space. Can we use less labelled data for training ?

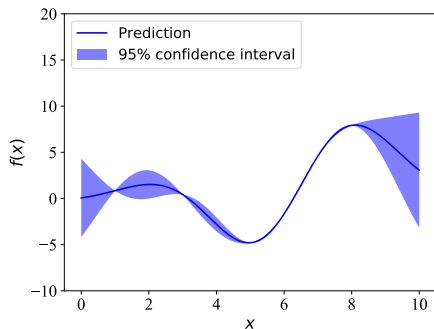
Solution: *Active Learning*

- Select representative samples to learn the property of the Pareto-Frontier;
- Fewer samples are needed, and the cost for labeling data is reduced.



Regression

- Gaussian process model;
- A prediction consists of a mean and a variance;
- Off-the-shelf library for implementation.



²Y. Ma, S. Roy, J. Miao, J. Chen and B. Yu, "Cross-Layer Optimization for High Speed Adders: A Pareto Driven Machine Learning Approach", TCAD'19.

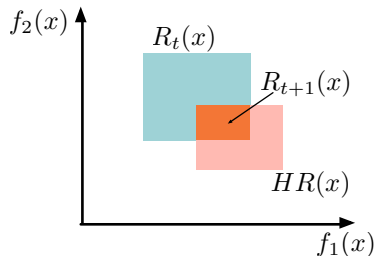
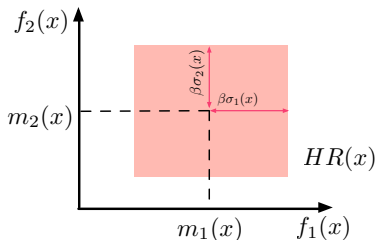
Uncertainty

- Given the prediction (m, σ) , a hyper-rectangle is defined as

$$HR(x) = \{y : m_i(x) - \beta\sigma_i(x) \leq y_i \leq m_i(x) + \beta\sigma_i(x)\}$$

- The uncertainty region is defined as:

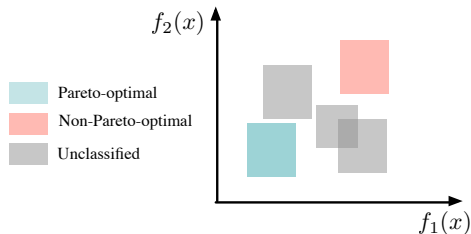
$$R_{t+1}(x) = R_t(x) \cap HR(x)$$





Classification

$$x \in \begin{cases} P, & \text{if } \max(R_t(x)) \leq \min(R_t(x')) + \delta, \\ N, & \text{if } \max(R_t(x')) \leq \min(R_t(x)) + \delta, \\ U, & \text{otherwise.} \end{cases}$$

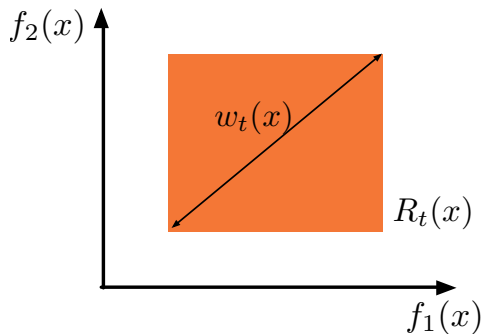


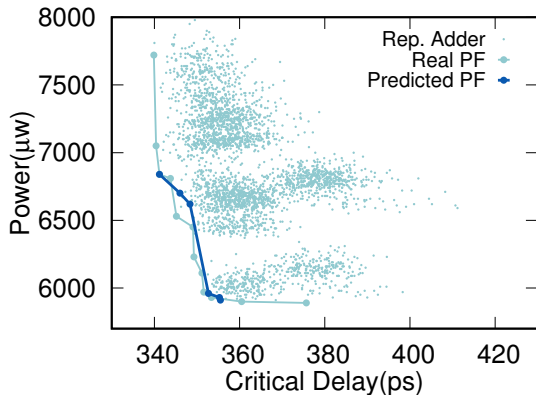
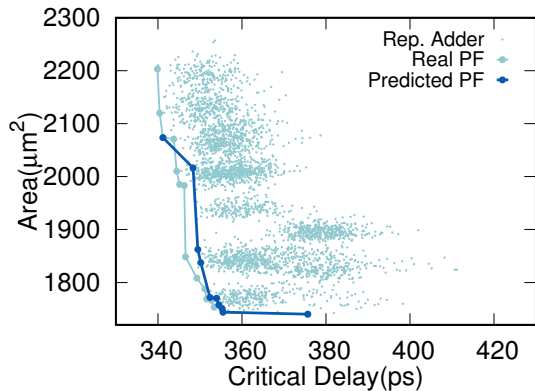


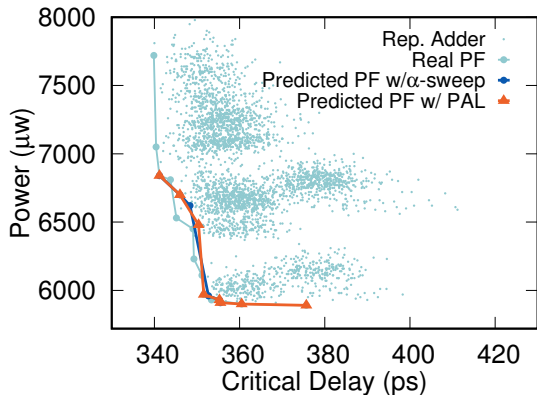
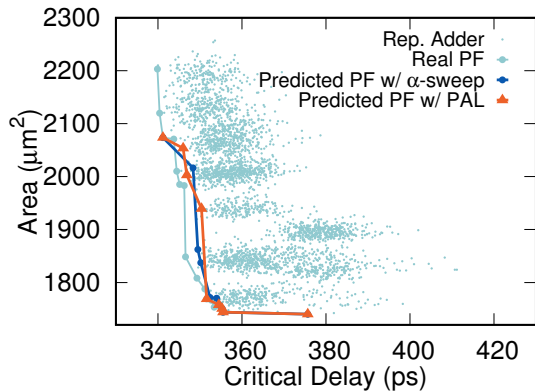
Sampling

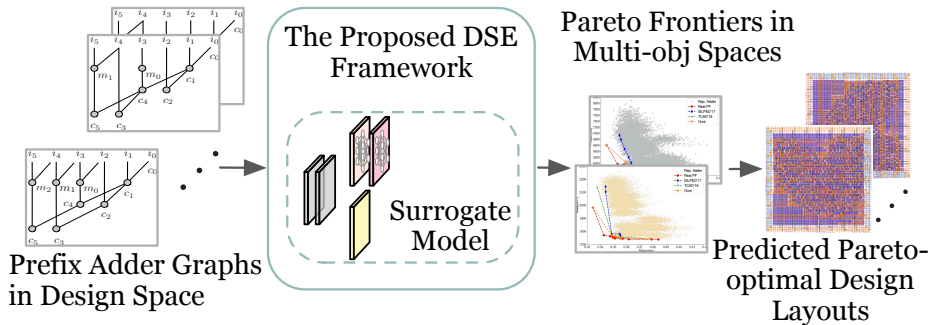
- Pick the one with largest uncertainty among the Pareto-optimal designs and unclassified designs.

$$w_t(x) = \max_{\mathbf{y}, \mathbf{y}' \in R_t(x)} \|\mathbf{y} - \mathbf{y}'\|_2.$$







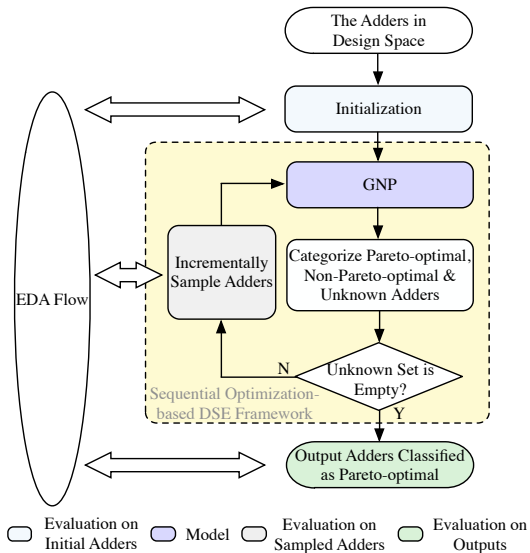


Graph Neural Process in the adder design space exploration.

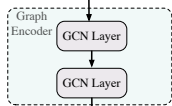
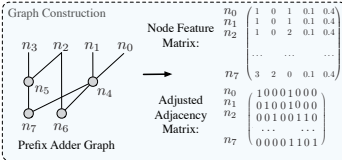
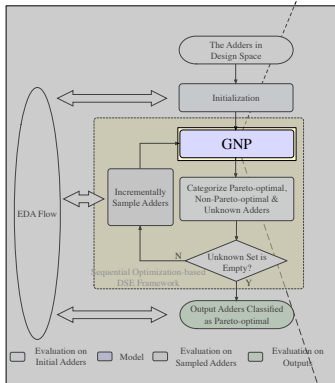
- A variational graph autoencoder is built to extract features from prefix adder structures automatically;
- A neural process is exploited to reduce computational complexity.

³H. Geng, Y. Ma, Q. Xu, J. Miao, S. Roy and B. Yu, "High-Speed Adder Design Space Exploration via Graph Neural Processes," TCAD'21.

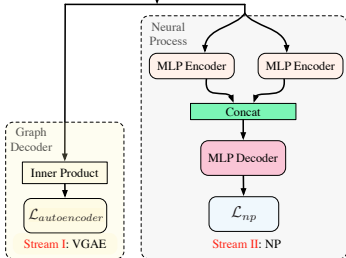
Our solution: Overview (I) [TCAD'21]



Our solution: Overview (II) [TCAD'21]



Latent Feature Vector of Adder Design

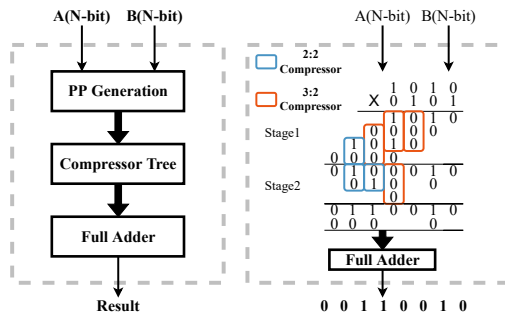


Output: Reconstructed Adjusted Adjacency Matrix

Output: Predicted QoR metric values with Uncertainties



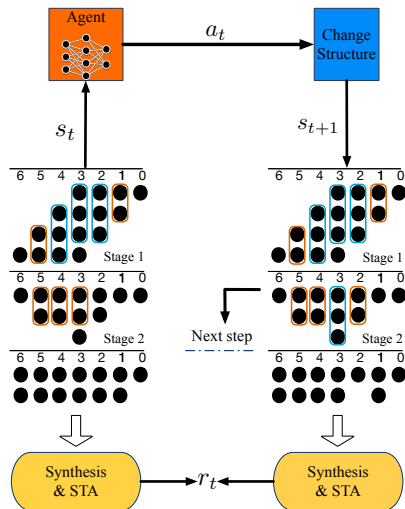
- A partial product generator (PPG)
- A compressor tree (CT), which is the most critical part.
- A carry propagation adder



Multiplier architecture

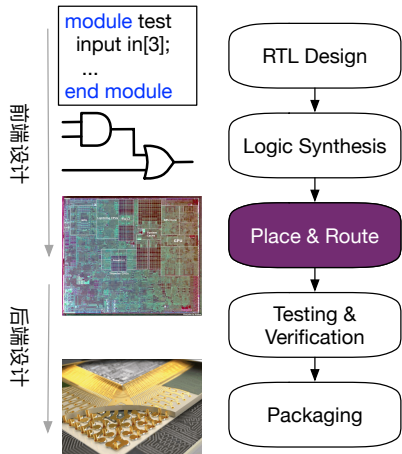
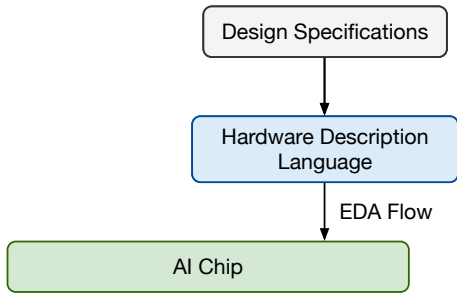
Deep Q-learning(DQN) based framework

- ResNet-18 as the agent
- A state s refers to a structure.
- An action a refers to modification on current structure s
- Pareto-driven Reward

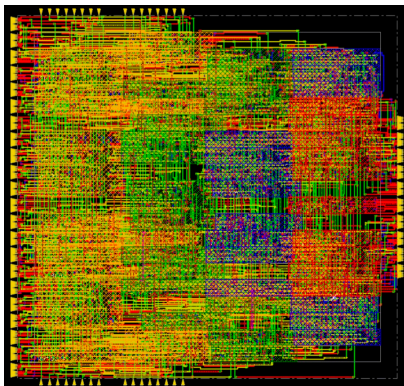
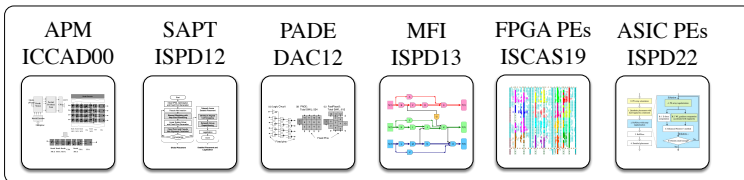


⁴D. Zuo, Y. Ouyang, Y. Ma "RL-MUL: Multiplier Design Optimization with Deep Reinforcement Learning," DAC'23.

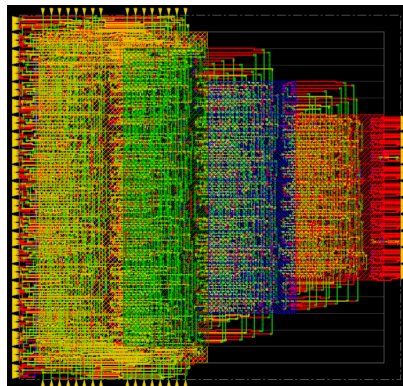
Datapath Driven Placement



Datapath
Placement



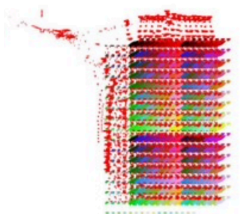
(a)



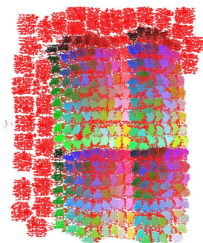
(b)



- Dataflow optimization becomes essential for AI chips
 - Schedules operation by data availability
 - Exposes opportunity for parallelism and data reuse
- Datapath regularity gives rise to new physical synthesis approaches
 - Directly determines system performance!



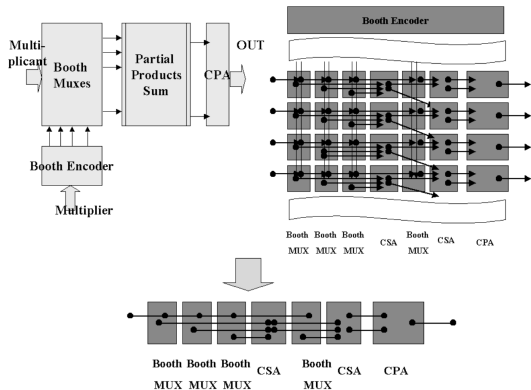
(a)Our flow 200 iteration



(b)Our flow detailed placement



- Datapath driven physical synthesis is **not new!**
- Placer has little control of exact locations if datapath is generated separately
- Abstract physical model
 - Bit-sliced abstraction
 - Compiled from HDL
 - Blocks are placed abutted
- Two-step heuristic for linear placement
 - quadratic assignment
 - sliding window optimization



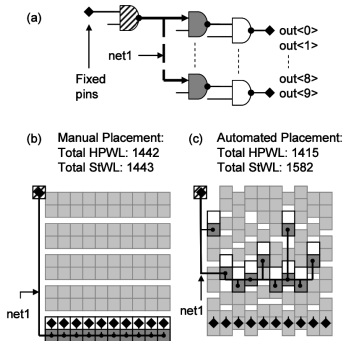
APM of a booth multiplier

¹Ye, T. Tao, and Giovanni De Micheli. "Data path placement with regularity", ICCAD'00.

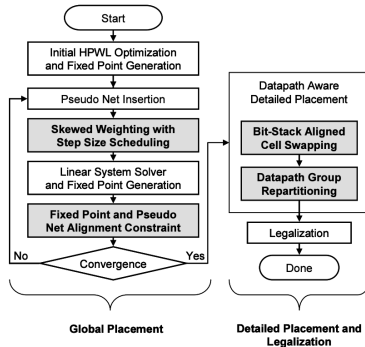


Datapath Driven Placement: SAPT [ISPD'12]²

- Structure aware placement obtains better steiner wire length → better routed wire length
- Apply bit-slice alignment constraints for force-directed placer
- Regularity driven detailed placement



Regularity matters

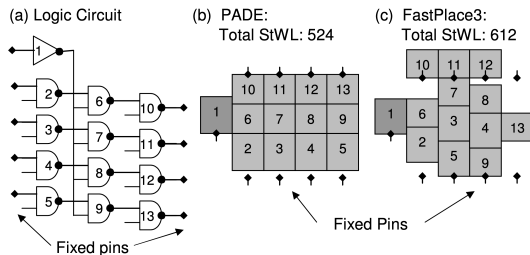


Structure aware placement

²Ward, Samuel I., et al. "Keep it straight: Teaching placement how to better handle designs with datapaths", ISPD'12.



- Machine learning techniques are involved
 - Graph-based (e.g., automorphism) and physical features (e.g., cell area) are analyzed and extracted from the netlist
 - Features are fed to SVMs and NNs to classify and evaluate datapath patterns
 - Maximize the evaluation accuracies of datapath and non-datapath pattern
 - Proposed new placer: PADE

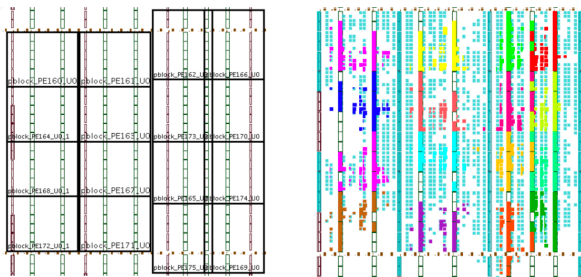


PADE effectively handles datapath in placement.

³Ward, Samuel, Duo Ding, and David Z. Pan. "PADE: A high-performance placer with automatic datapath extraction and evaluation through high dimensional data learning", DAC'12.

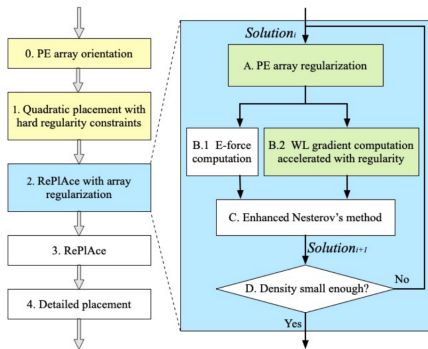


- Systolic arrays are a popular choice to support neural network computations
- Current FPGA CAD tools cannot synthesize them in high quality
- One solution: restrict fixed locations for PEs
 - Sufficient DSPs, close to used I/O banks

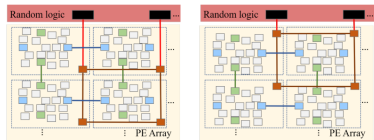


PE placement with floorplan constraints

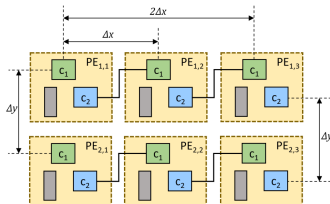
⁴Zhang, Jiayi, et al. "Frequency improvement of systolic array-based CNNs on FPGAs", ISCAS'19.



Placement Flow



(a) PE cells are placed in macro gen- (b) PE cells are placed along with random logic, wirelength: 244



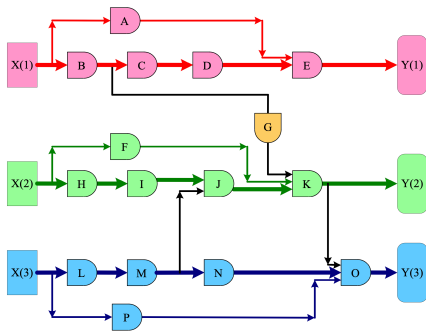
Regularity Constraints

- Simultaneously place the PEs and random logic cells for better solution quality
- Analytical global placement with PE array regularization

⁵Fang, Donghao, et al. "Global placement exploiting soft 2D regularity", ISPD'22.



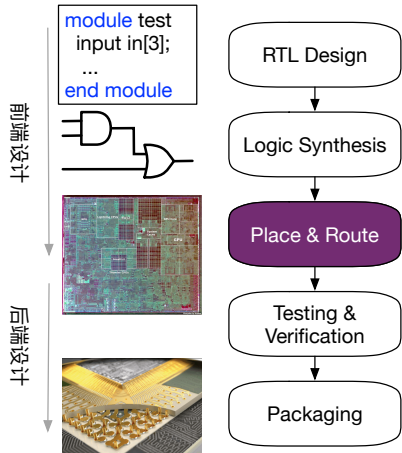
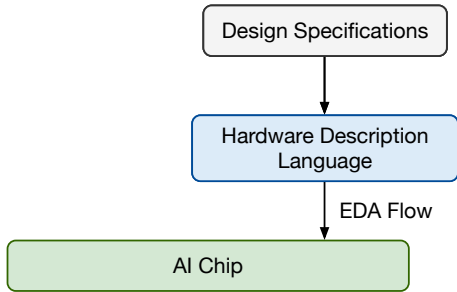
- *Datapath main frame (DMF)*
 - a set of n disjoint paths from input to output
 - maximize the number of datapath gates on these paths
- Can be optimally identified by the min-cost max-flow algorithm



DMF identification can be transformed to a network flow problem

⁶Xiang, Hua, et al. "Network flow based datapath bit slicing", ISPD'13.

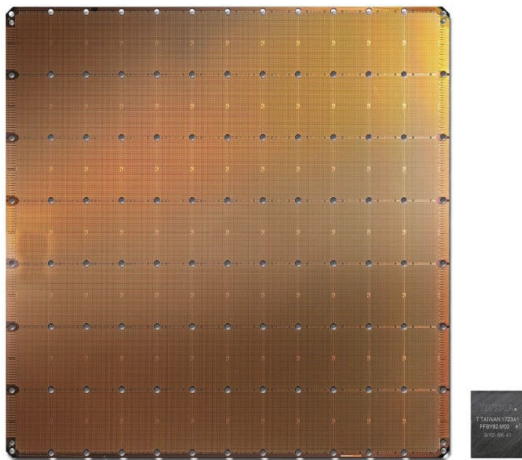
Wafer-Scale Floorplan





The Cerebras CS-1's Wafer-Scale Engine (WSE) is the largest and most powerful processor ever built.

- Consisting of 800 by 1060 identical processing elements.
- Keep all memory and all computation together on a single monolithic chip.

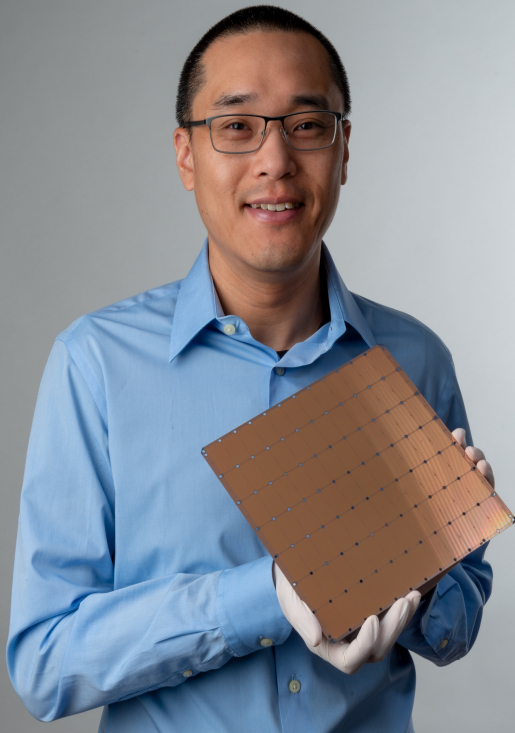


Die photograph of the Cerebras Wafer-Scale Engine (WSE; at left). For comparison, the largest GPU is shown on the right to scale.

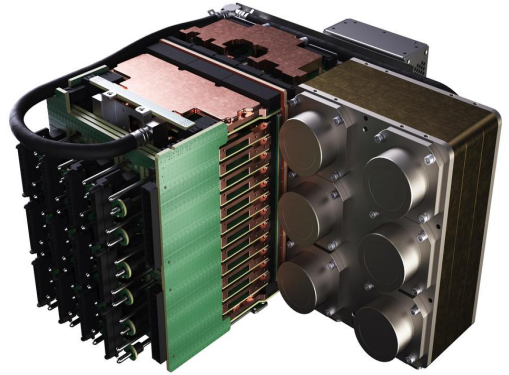
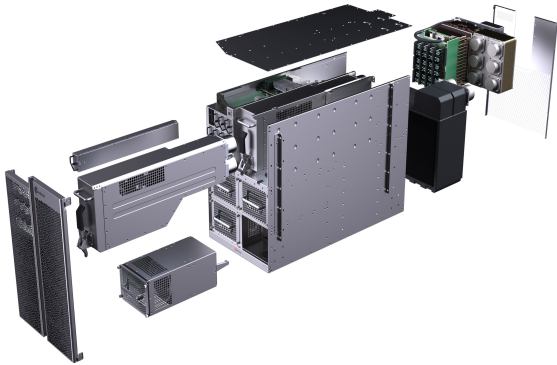
⁷J.P. Fricker, A. Hock, "Building a Wafer-Scale Deep Learning System: Lessons Learned", Supercomputing'19.

Cerebras: Largest Chip Ever Built

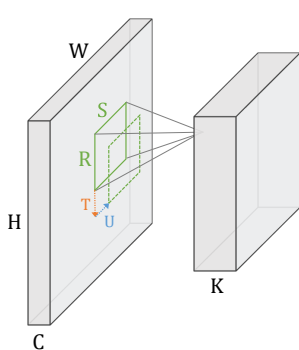
- 46,225 mm² silicon
- 2.6 **Trillion** transistors, 7nm TSMC
- 850,000 processor cores
- 40 Gigabytes of On-chip Memory
- 9 PByte/s memory bandwidth
- 100 Pbit/s fabric bandwidth



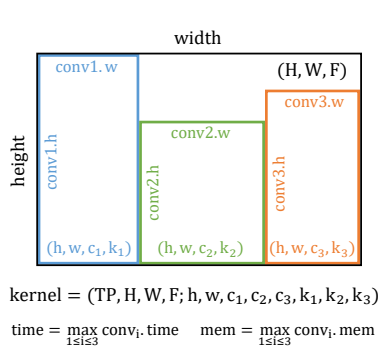
CS-1 Supercomputer Hardware



- conv*: basic convolution kernel



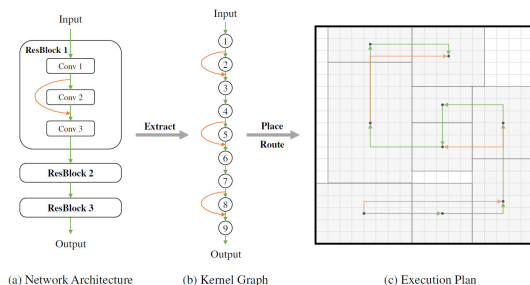
(a) Arguments of *conv*



(b) Performance of a kernel with 3 *conv*s

- 8 formal arguments: (H, W, R, S, C, K, T, U) \Rightarrow fixed input parameters.
- 4 execution arguments: (h, w, c, k) \Rightarrow variables to be determined.

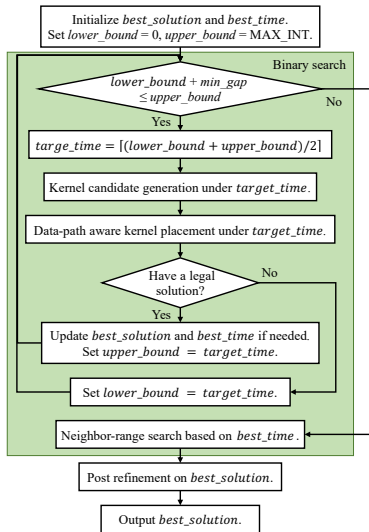
- Determine the **execution parameters** and the **locations** for all kernels. [ISPD'20]⁸
 - Blocks are soft: kernel sizing
 - Floorplanning
 - Optimize performance, wirelength, etc.



Kernel floorplanning. Figure adopted from [ICCAD'20]⁹

⁸Michael James, et.al., "Physical Mapping of Neural Networks on a Wafer-Scale Deep Learning Accelerator", ISPD'20.

⁹Bentian Jiang, et.al., CU.POKer: Placing DNNs on Wafer-Scale AI Accelerator with Optimal Kernel Sizing, ICCAD'20.



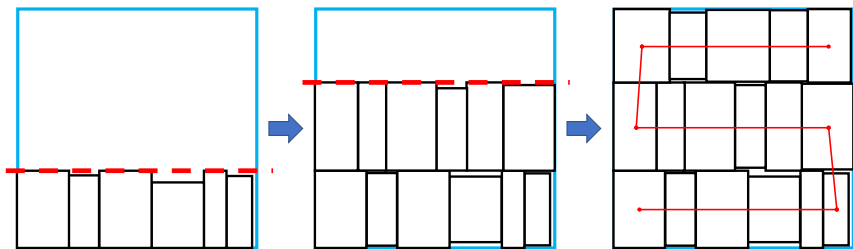
Two-steps Search

- Binary search
- Neighbor-range search
- Post refinement

Searching under Target Time

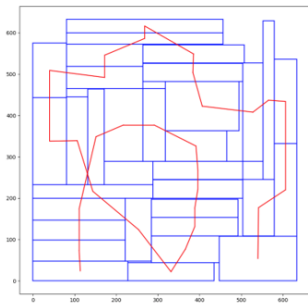
- Kernel candidates generation with optimal shapes **under given target time**
- **Data-path aware** placement

⁹Bentian Jiang, et.al., CU.POKer: Placing DNNs on Wafer-Scale AI Accelerator with Optimal Kernel Sizing, ICCAD'20.



Overall Flow

- Given a target time T , generate all the kernel candidates with optimal shapes and execution times under T .
- According to the connectivity graph, generate the topological order of the kernels for placement.
- Place the kernels compactly row by row in the topological order.



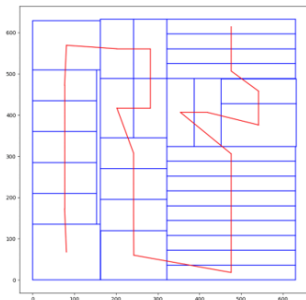
Simulated Annealing Placer:

Max_time: 76698

Wire_length: 3237

Adapter_cost: 15

Score: 110478



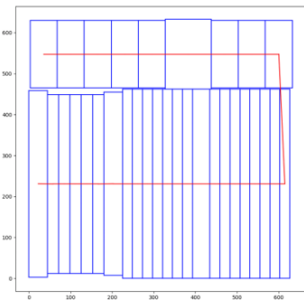
Divide and Conquer Placer:

Max_time: 65106

Wire_length: 2650.5

Adapter_cost: 18

Score: 93321



CU.POKer:

Max_time: 65170

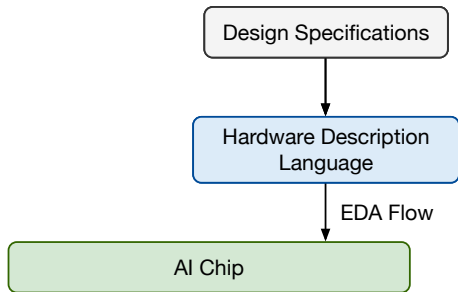
Wire_length: 1489.5

Adapter_cost: 12

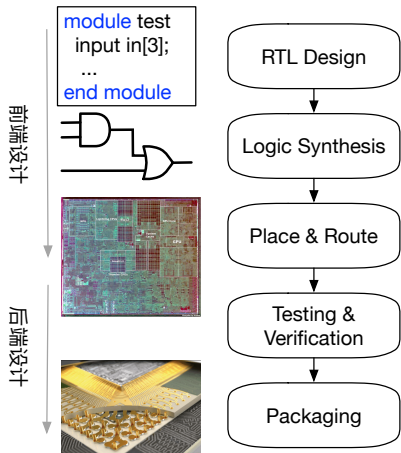
Score: 81265

Performance comparisons with SA and DC placers on 8 public benchmarks.

Conclusion



- Arithmetic Unit Synthesis
- Datapath Driven Placement
- Wafer-Scale Floorplan



THANK YOU!