

An Optimization Playground for Precision and Number Representation Tuning

Olivier Sentieys

Univ. Rennes, Inria

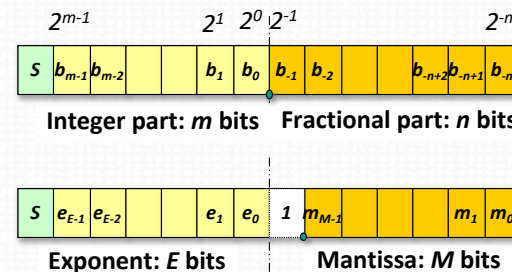
with Van-Phu Ha, Tomofumi Yuki, Daniel Ménard and others

Inria



Computer Arithmetic

- At the core of computing we find **number representations** (e.g., real) and **arithmetic operations** (e.g., +, \times , \div , $\sqrt{\quad}$)
- Numeric formats
 - Fixed-Point (integer)
 - Floating-Point
- Energy, delay, and area vary a lot between numeric formats and word-length

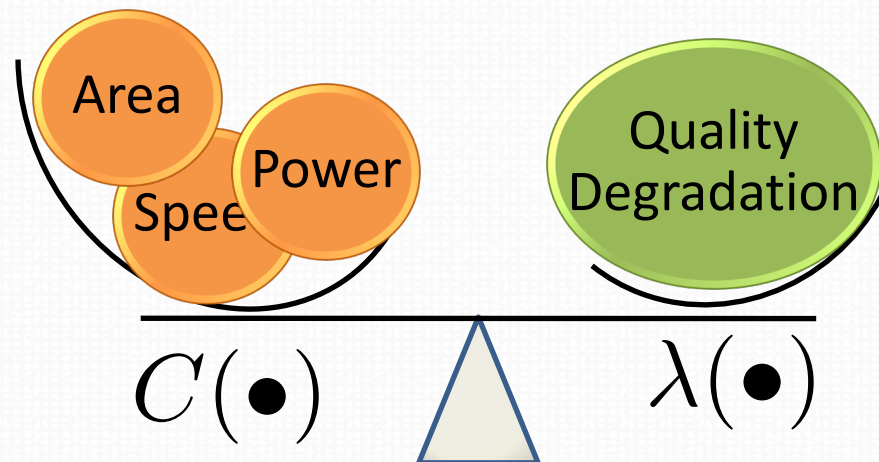


	Addition	Multiplication
8-bit integer	0.03pJ / 36 μm^2	0.2pJ / 282 μm^2
32-bit float	0.9pJ / 4184 μm^2	3.7pJ / 7700 μm^2

Automatic Precision Tuning

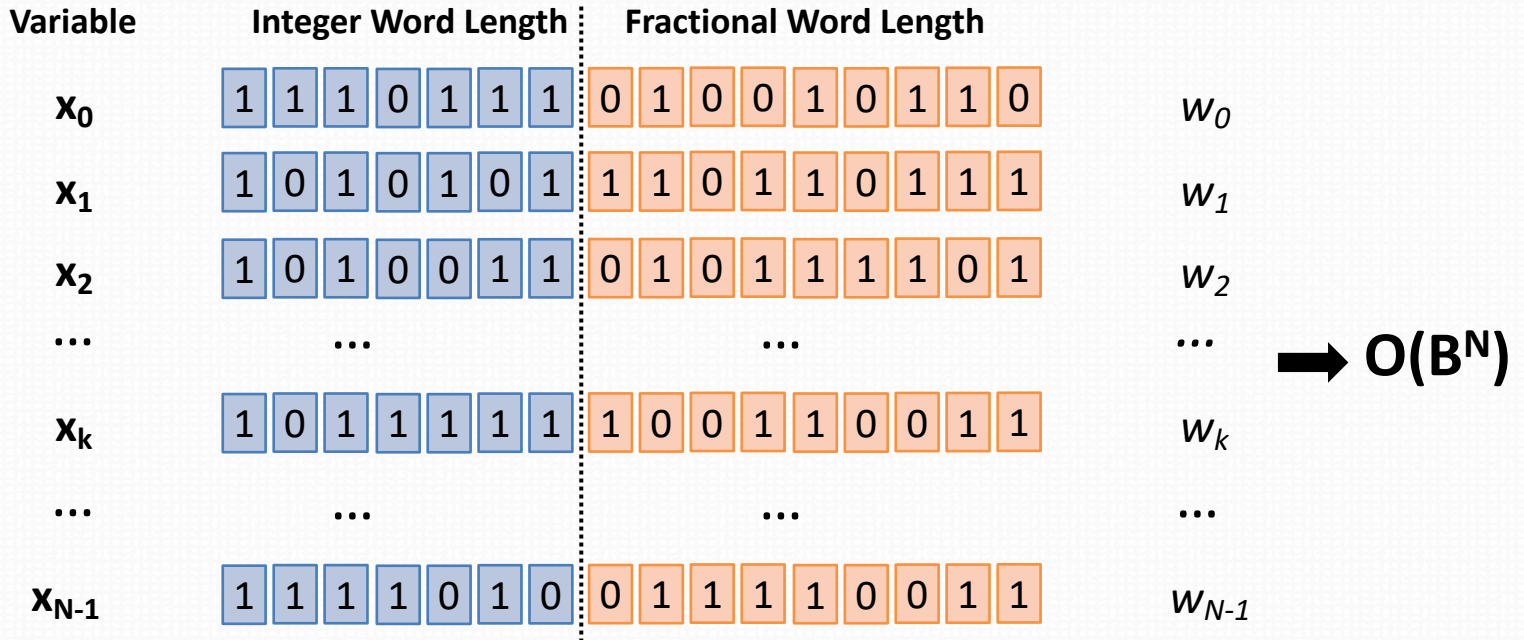
[also Word-Length Optimization (WLO)]

- Optimization process that
 - determines the number of bits for each data
 - minimizing a cost function C
 - constrained by (application) quality degradation λ
 - e.g., noise power, SSIM, abs. error



Automatic Precision Tuning

Multiple Word Length



N: number of variables

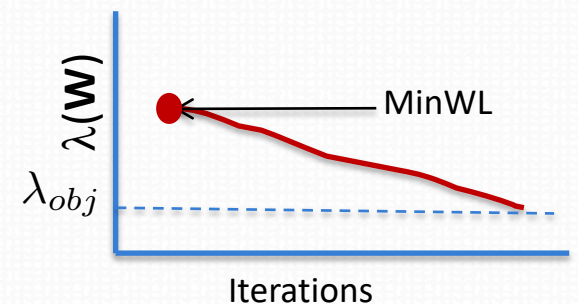
B: number of bits to explore per variable

Automatic Precision Tuning

- Multi-variable **word-length optimization**

$$\min (C(\mathbf{w})) \quad \text{subject to} \quad \lambda(\mathbf{w}) \leq \lambda_{obj}$$

- Known to be **non-convex** and **NP-hard**
- Optimized using heuristic rules, iterative optimization process, stochastic approaches



$\lambda(\mathbf{w})$: accuracy degradation of solution \mathbf{w}

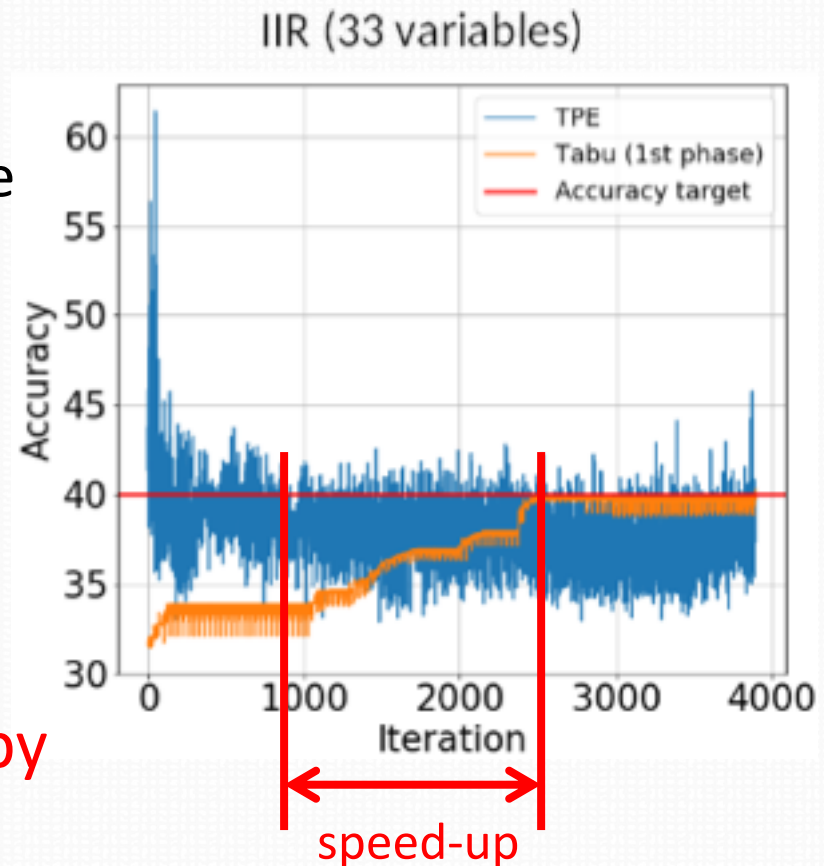
$C(\mathbf{w})$: cost of solution \mathbf{w}

Data word lengths: $\mathbf{w} = \{w_0, w_1, \dots, w_{N-1}\}$

Maximum degradation: λ_{obj}

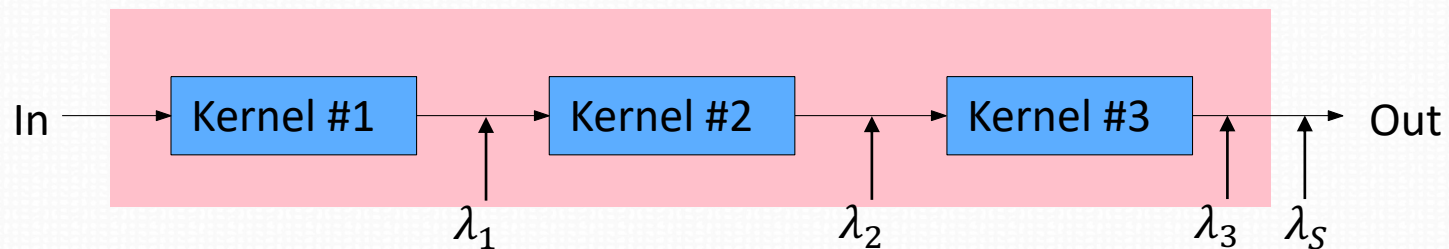
Speeding-up Global Search

- Combine **Bayesian Optimization** and **Local Search**
 - Bayesian Optimization for narrowing down solution space
 - Fine-tuning with local search
- Transition point based on statistical metrics
 - word-lengths (WL) are distributed with low variance
 - e.g., with less than 1 bit
- **Optimization time is reduced by 50-80%** w.r.t. best algorithm with similar cost



Scaling the WLO Procedure

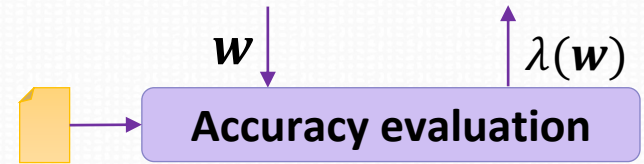
- Large system sizes present enormous complexity
 - Too many variables for global optimization



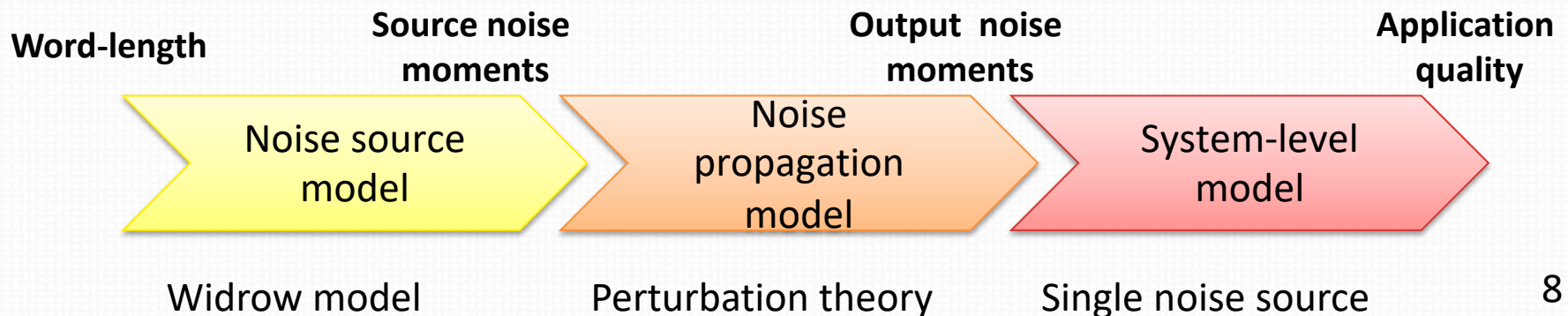
Multi-kernel approach

- Key idea: **construct models** that express
 - impact of **noise budgets** to Cost and Accuracy
 - relation among **noise budgets**
- Significantly **reduce exploration time** and improve the quality of the solutions for large applications

Accuracy Evaluation

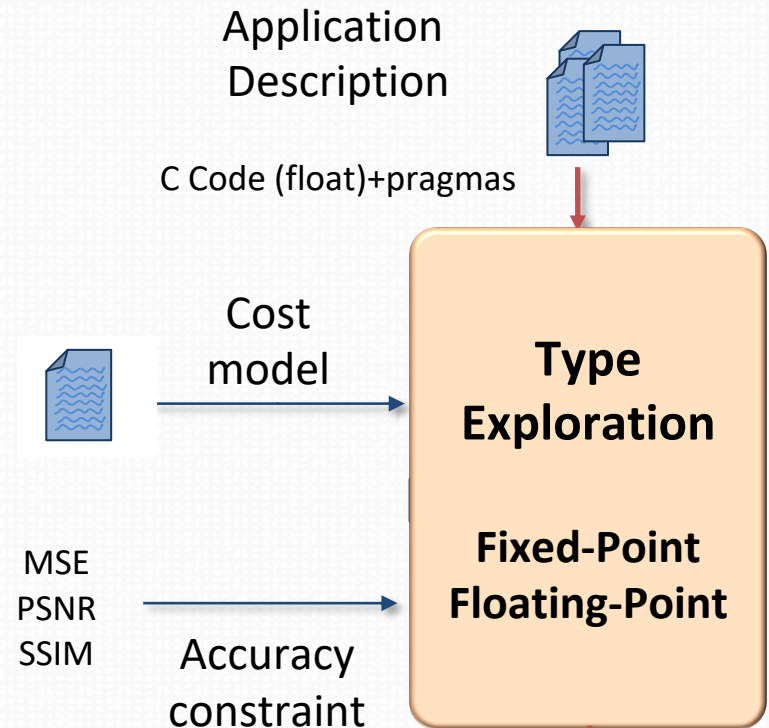
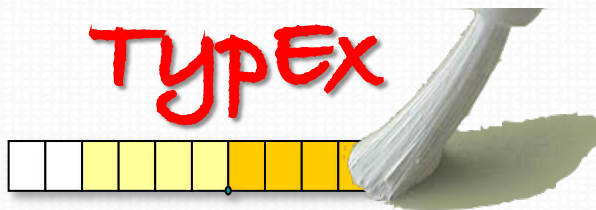


- One of the most time consuming tasks during precision tuning
- Models for quantization effect analysis
 - **Analytical** accuracy evaluation
 - **System-level** estimation [ICCAD'14, DATE'16]
 - **Speeding-up** simulations [DATE'20, ICCAD'14]



TypEx: A Framework for Type Exploration

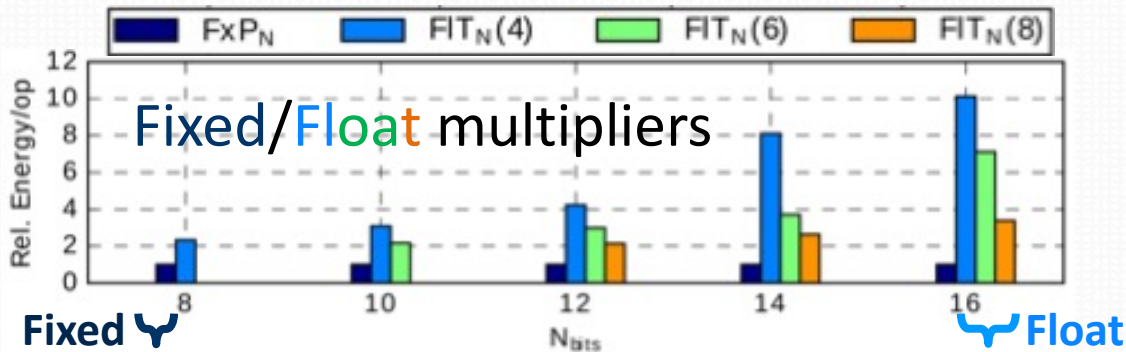
- Source-to-source
 - C code in float to C code using custom arithmetic
- Word-length optimisation
 - fixed or float



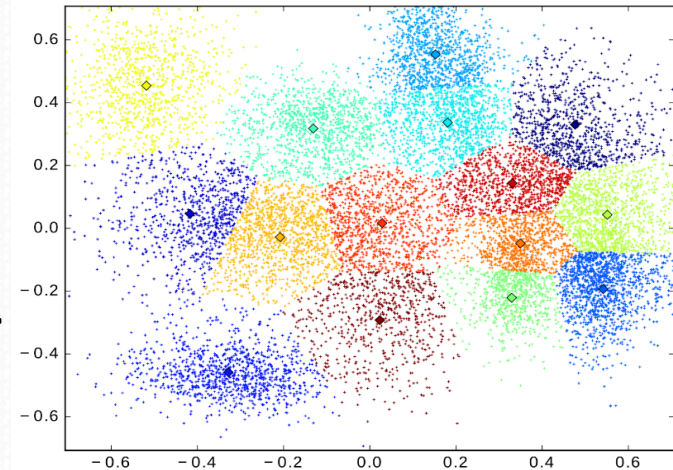
C++ Code
Customized Arithmetic

Custom Floating-Point

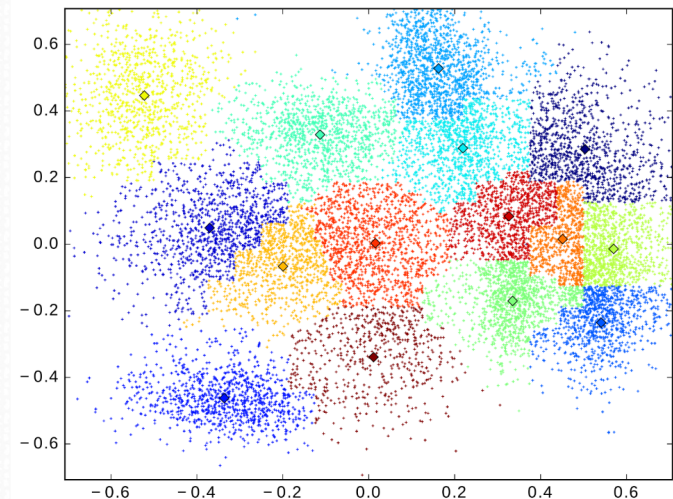
- Slower increase of errors for floating-point
 - e.g., 8-bit float is still effective for K-means clustering [SiPS'17]
- Difference in cost/energy between float/fixed is small for low-precision operators



Approximate K-Means Clustering



Reference: double



Floating-Point: ct_float₈

5-bit exponent

3-bit mantissa

10

Custom Floating-Point

- `ct_float`: a Custom Floating-Point

C++ Library

<https://gitlab.inria.fr/sentieys/ctfloat>

- **Synthesizable** (with HLS) library

- Templated C++ class

```
ct_float<e,m,r>
```

- Exponent width e (int)
- Mantissa width m (int)
- Rounding method r
- Bias b

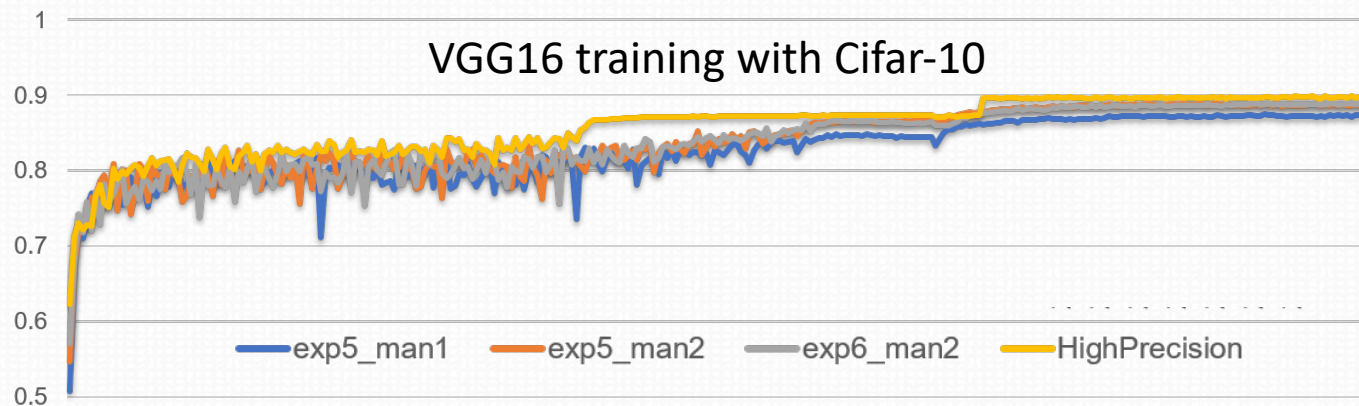
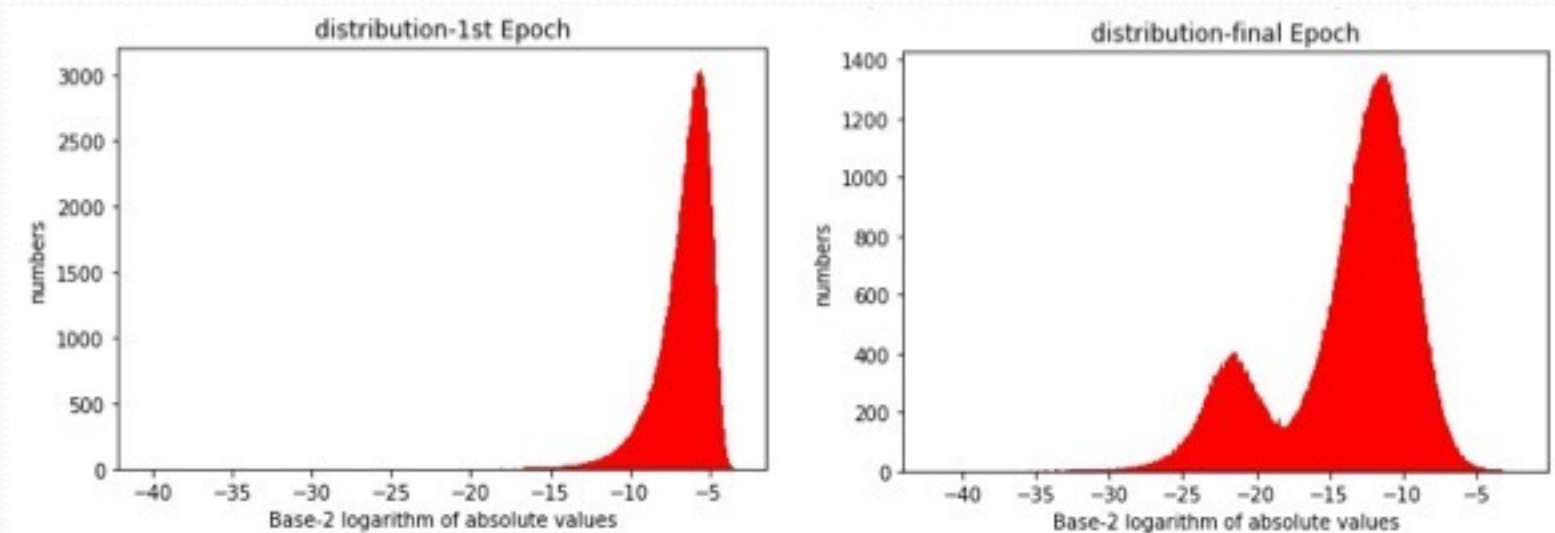
```
ct_float<8,12,CT_RD> x,y,z;  
x = 1.5565e-2;  
z = x + y;
```

- **Many possible design points**

- latency constraints, rounding modes, etc.

Work in Progress

- Low-Precision **Training** of DNNs



Open Issues

- WLO is still a difficult problem for **large** applications
 - Mainly limited by simulation time to evaluate $\lambda(\mathbf{w})$, analytical evaluation still limited
- Evaluating **cost** $C(\mathbf{w})$ is also an important (and less studied) issue
 - Resource sharing, complexity related to one w_i
- Automating the choice between (or combining) **float and fixed** is a challenge
 - Towards an automatic optimizing compiler framework
 - considering both float and fixed representations
 - combining WL/Number Representation optim.

More Open Issues

- Compiler level: identify candidate computation kernels for approximations
- Hardware-level: build a precision-reconfigurable acceleration platform, especially suitable for training ML/DL
- Algorithmic-level: build analytical (maths) models of accuracy loss in DNN to avoid long simulations