



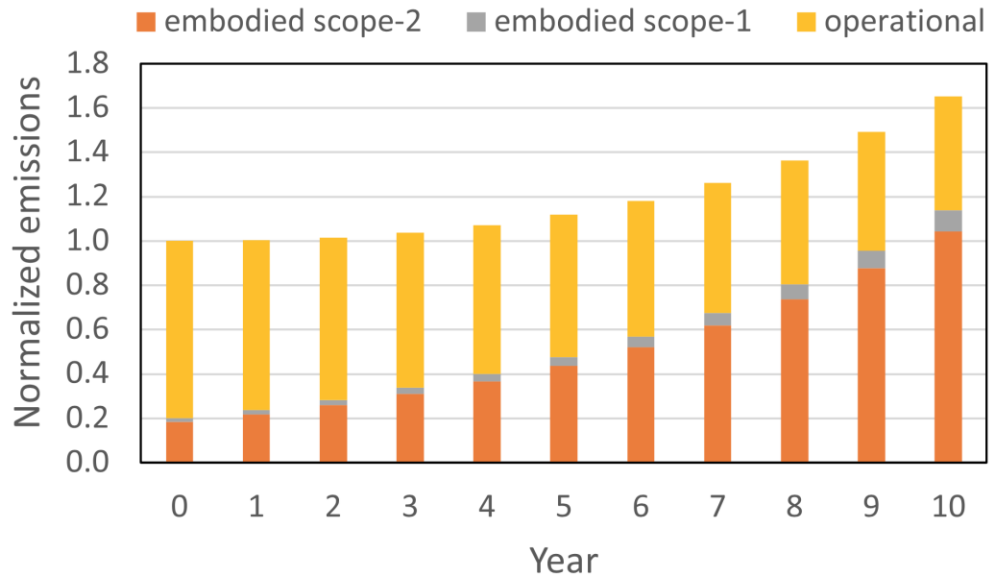
Wordlength Optimization for Custom Floating-point Systems

DASIP 2024

Authors :
Quentin Milot, Mickaël Dardaillon,
Justine Bonnot and Daniel Menard



Introduction



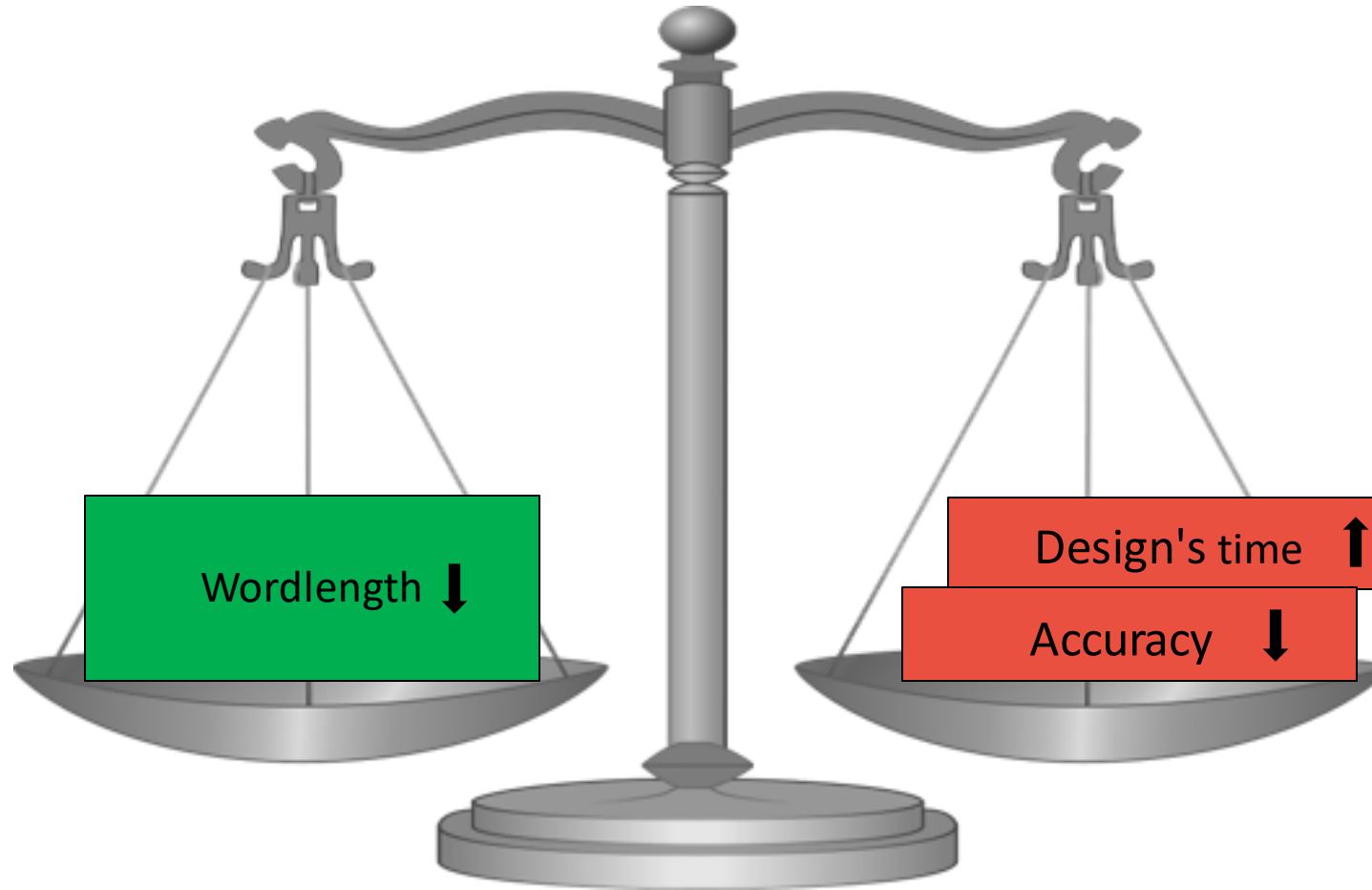
Extracted from [1].

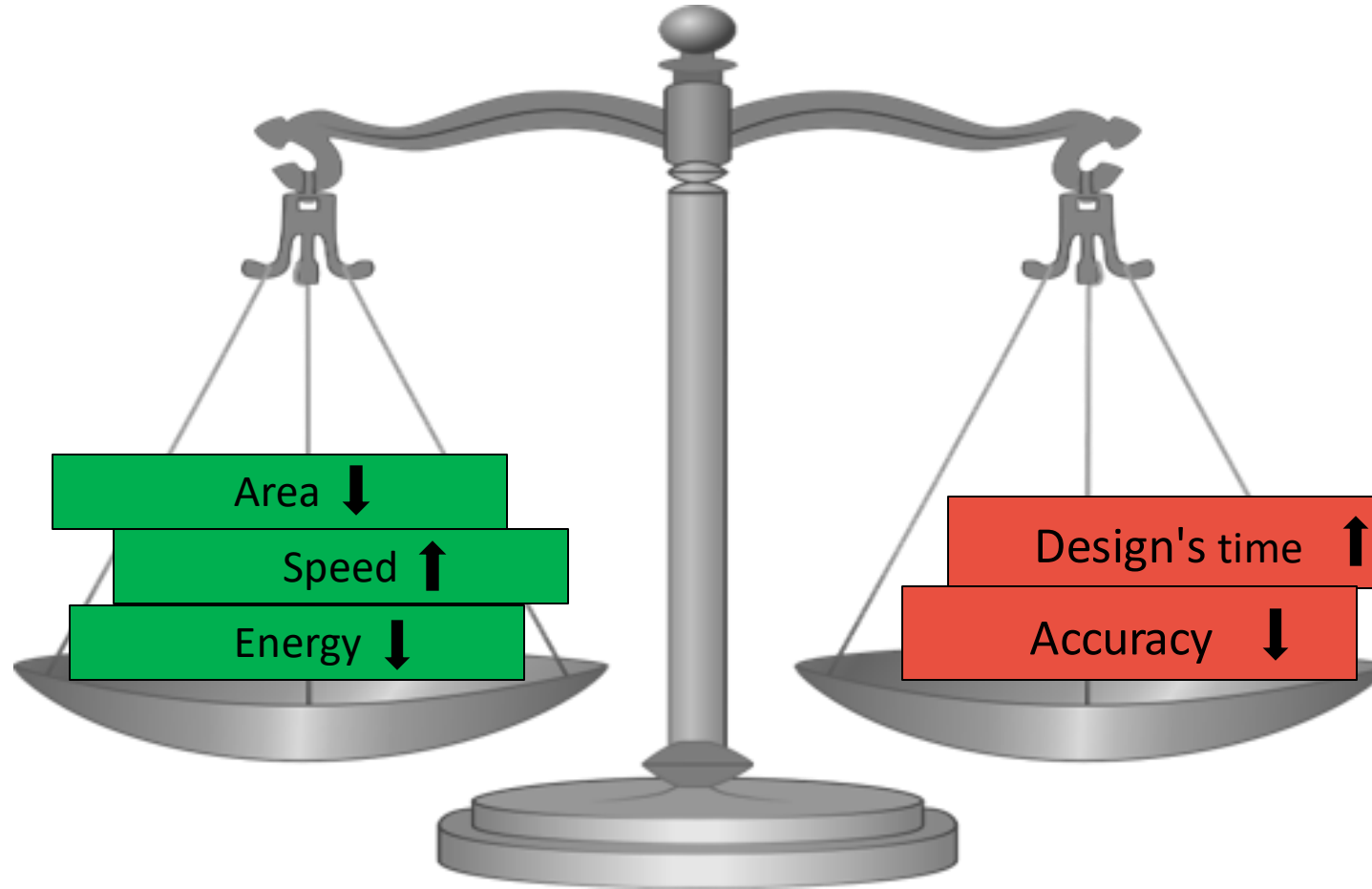
A scenario of normalized emissions projections for computer systems over the year

- Environmental impact of computer systems increase and will increase over the years
- Impact come from production, deployment and usage
- Proposed solution :
 - Reduce die size
 - Reduce operational energy

=> A lead : Use of approximate computing

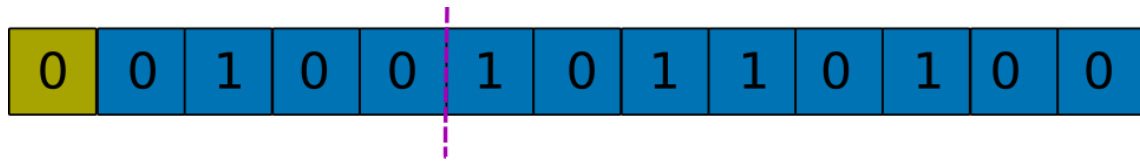
[1]"Kaya for Computer Architects: Toward Sustainable Computer Systems", Lieven Eeckhout, IEEE Micro, Volume 43, 2023



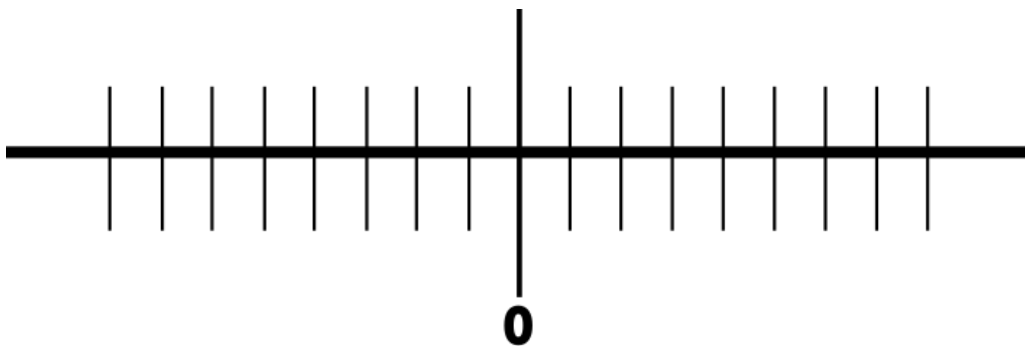


Fixed point

$$X * 2^{-E} = (-1)^S * Int$$



- Exponent implicit in the code

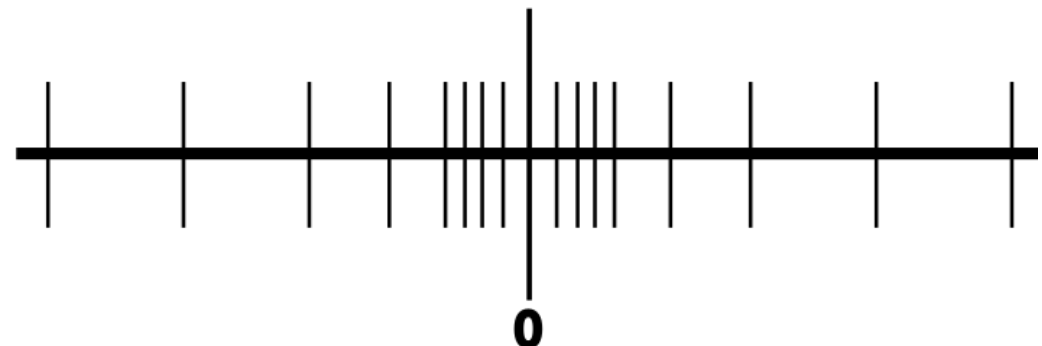


Floating point

$$X = (-1)^S * (1 + M) * 2^{E-\Delta}$$



- Exponent explicit in the coded value



Operation		Picojoules per operation	
		45nm	7nm
+	Int 32	0.1	0.03
	IEEE FP 32	0.9	0.38
X	Int 32	3.1	1.48
	IEEE FP 32	3.7	1.31
SRAM	8KB	10	7.5
	32KB	20	8.5
	1MB	100	14

Data extracted from [2]
Energy per operation for a 45nm and 7 nm systems

- Fixed point was first choice thanks to low operational energy
- Floating point operational energy reduces over the year (~ - 2X)
- Memory usage is the higher energy requirement now
- Floating point can lead to lower memory usage

=> Custom floating-point systems became an option for embedded systems

[2]"Ten Lessons From Three Generations Shaped Google's TPuv4i", Norman P. Jouppi et al, ISCA 2021

Challenge :

Determine a generic **automatic** wordlength optimization method for a given accuracy

Outline :

I. Introduction

II. Optimization flow

III. Proposed strategies

IV. Experiments

Outline :

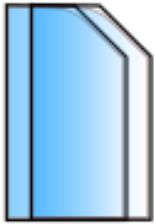
I. Introduction

II. Optimization flow

III. Proposed strategies

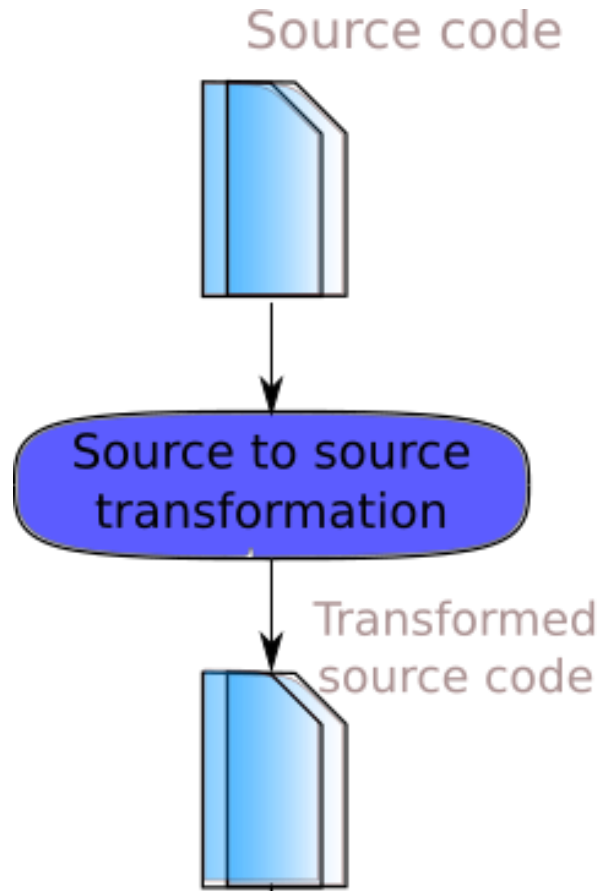
IV. Experiments

Source code



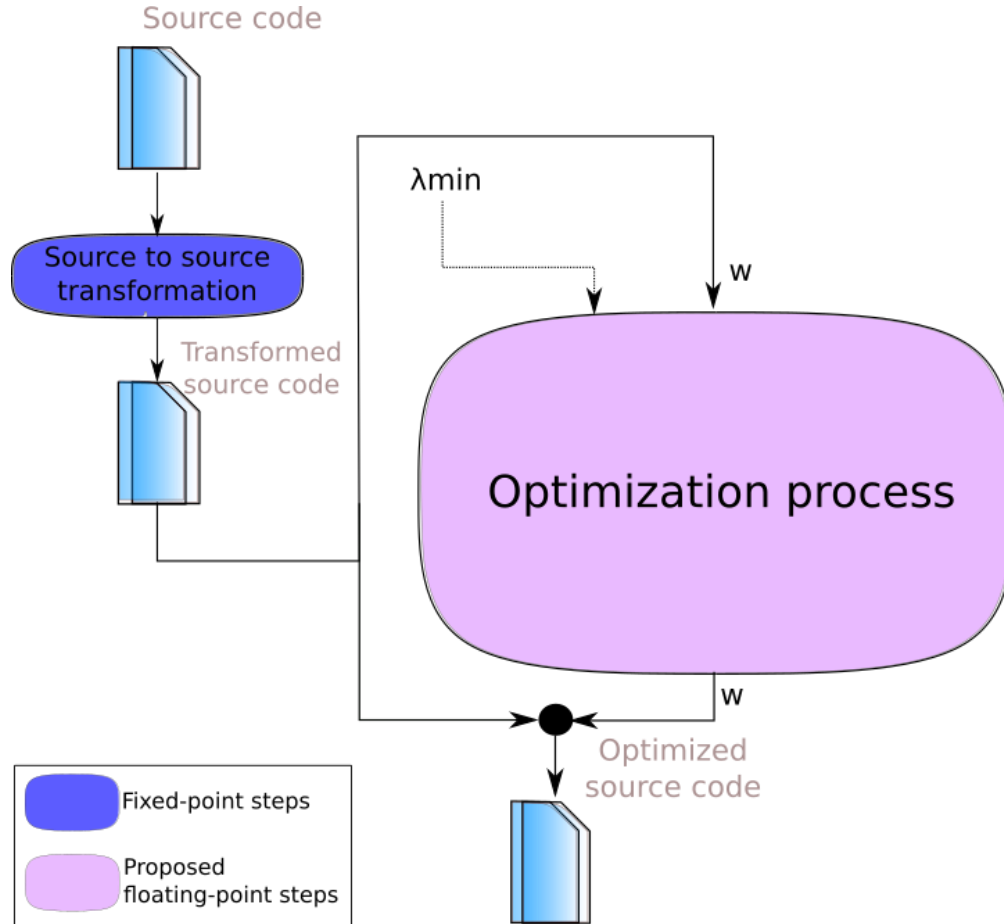
User source code :

- Write in C/C++
- Any algorithm



Source to source transformation :

- Write in C/C++
- Any algorithm



Optimization :

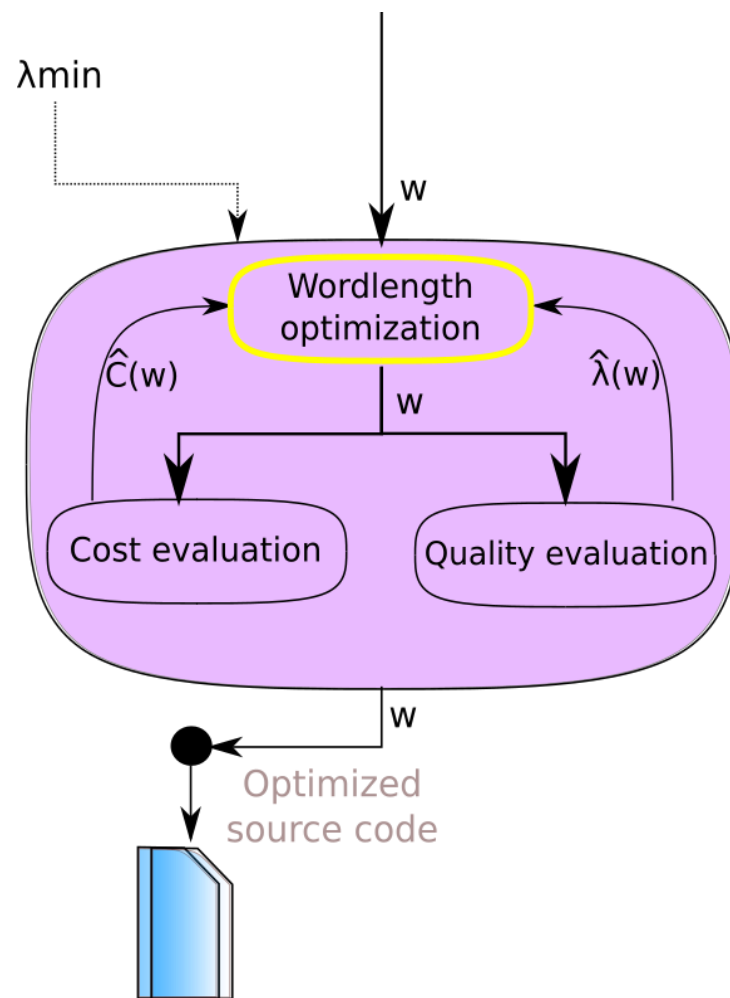
- Input : $(W_{E_{max}}, W_{E_{max}}, \dots, W_{E_{max}}, W_{M_{max}}, W_{M_{max}}, \dots, W_{M_{max}})$
- λ_{min} : The user minimal quality
- Problem :
$$\min_{\hat{\lambda}(W) \geq \lambda_{min}} \hat{C}(W)$$

Variable's exponent Variable's mantissa

Wordlength Optimization

MIN+1 algorithm:

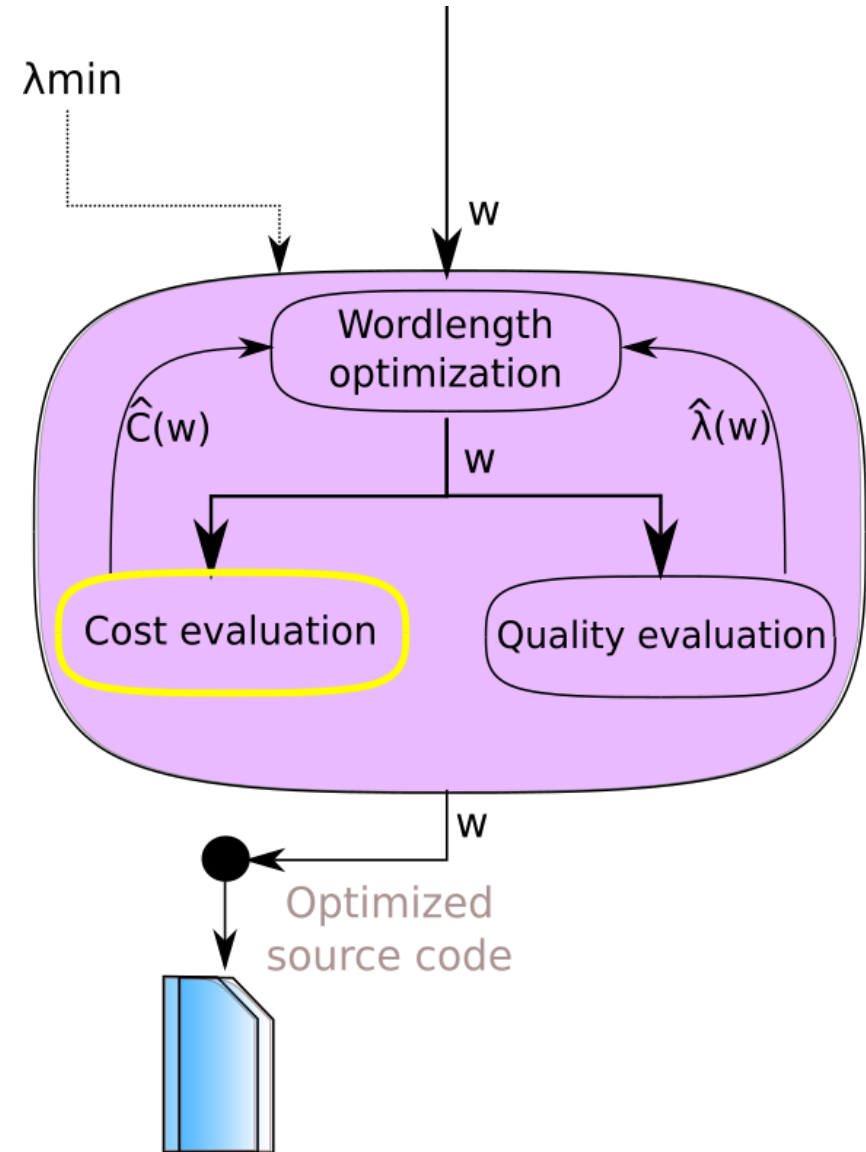
- Greedy heuristic algorithm
- Two steps algorithm :
 - 1) Find the minimal value for each dimension that satisfy the criterion
 - Set all value to the minimal found
 - 2) Increment the dimension with the best gradient until the criterion is met



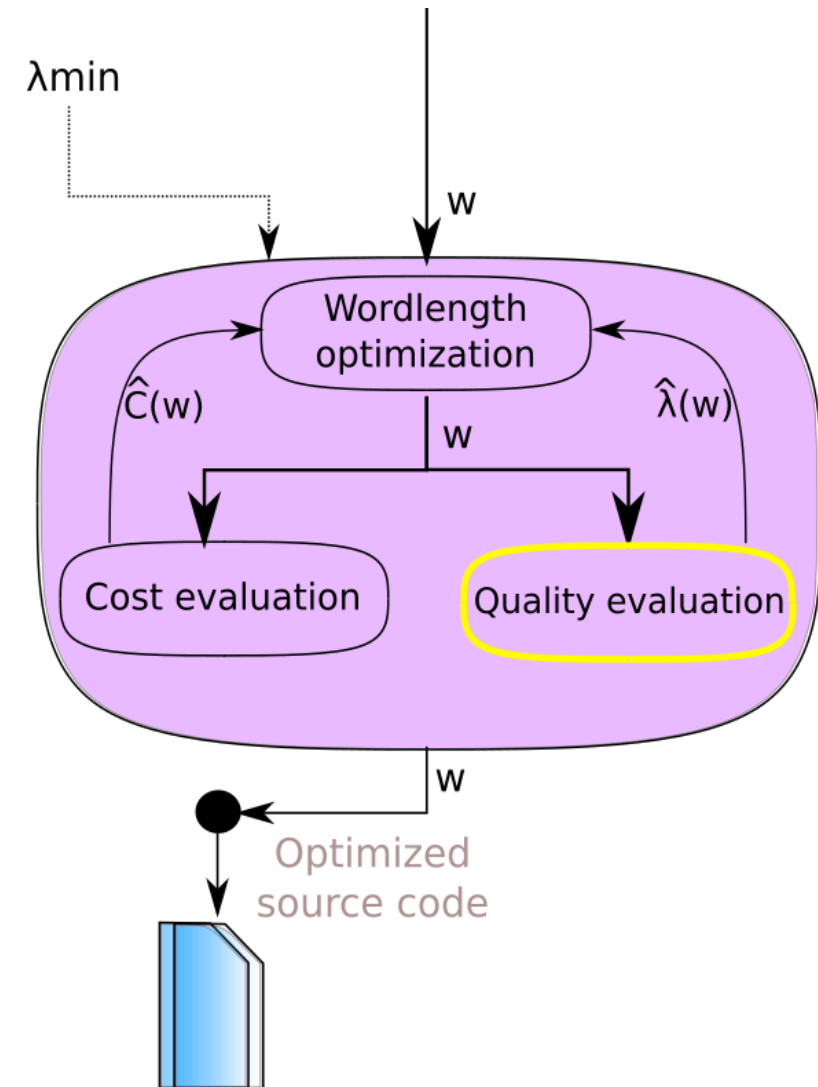
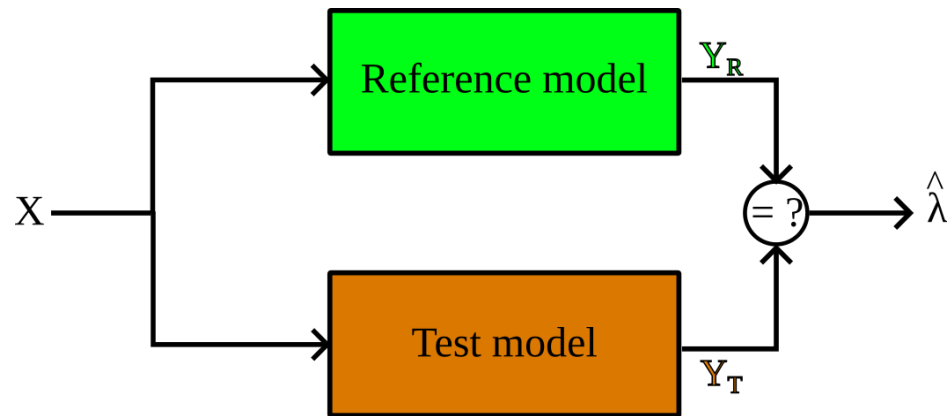
Cost evaluation

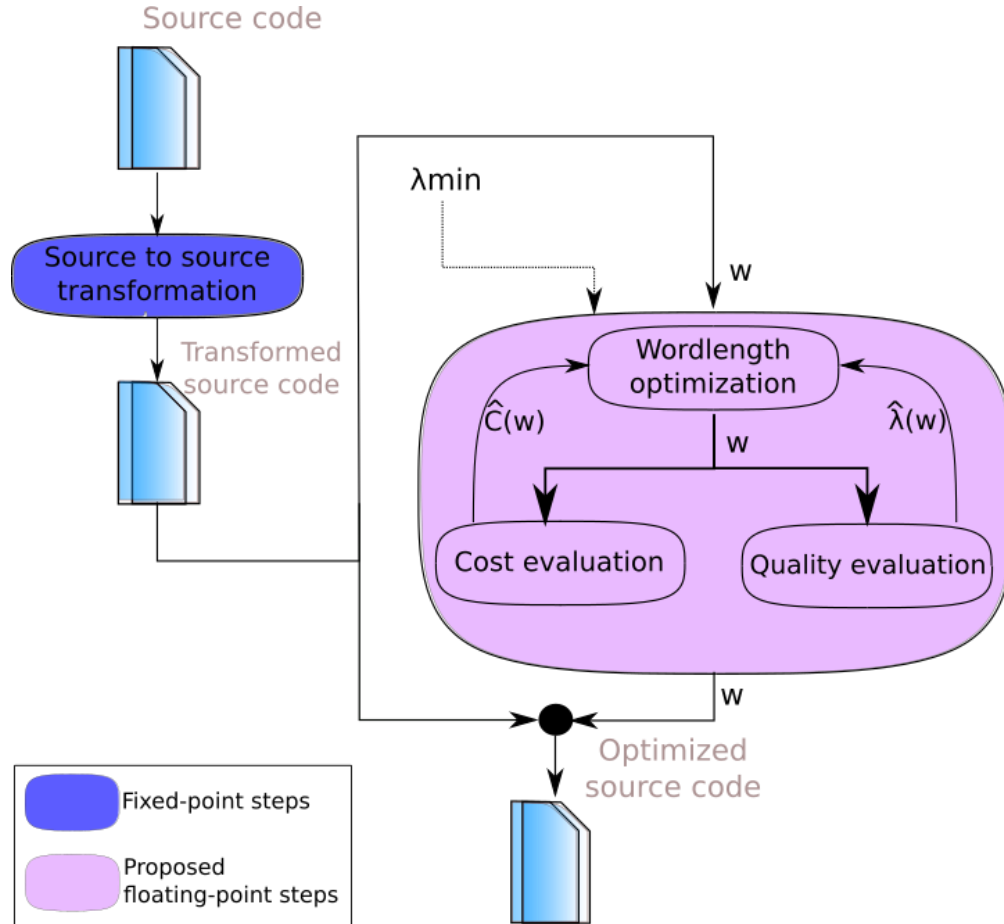
Requirement :

- Types of memory
- Energy for a given operation and wordlength
- Architectures



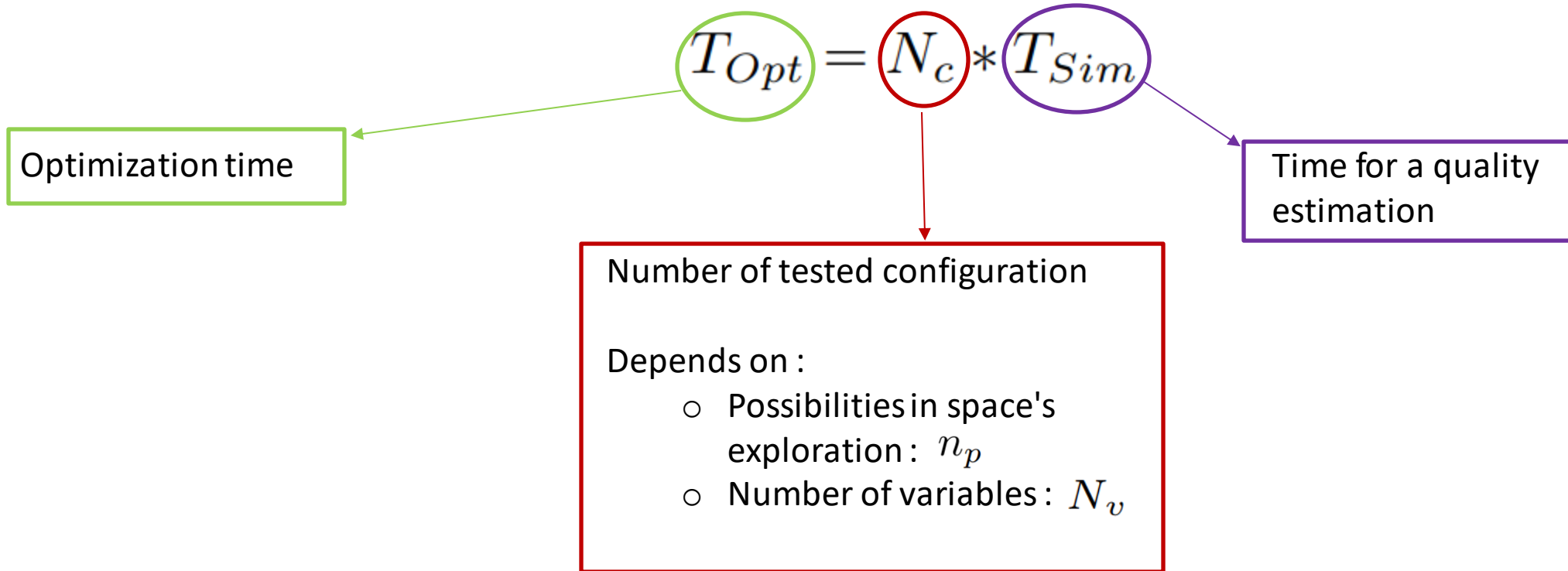
Quality evaluation





Optimization process :

- Input : $(W_{E_{max}}, W_{E_{max}}, \dots, W_{E_{max}}, W_{M_{max}}, W_{M_{max}}, \dots, W_{M_{max}})$
- λ_{min} : The user minimal quality



Challenge :

Determine a generic **automatic** wordlength optimization method for a given accuracy in a **reasonable amount of time**

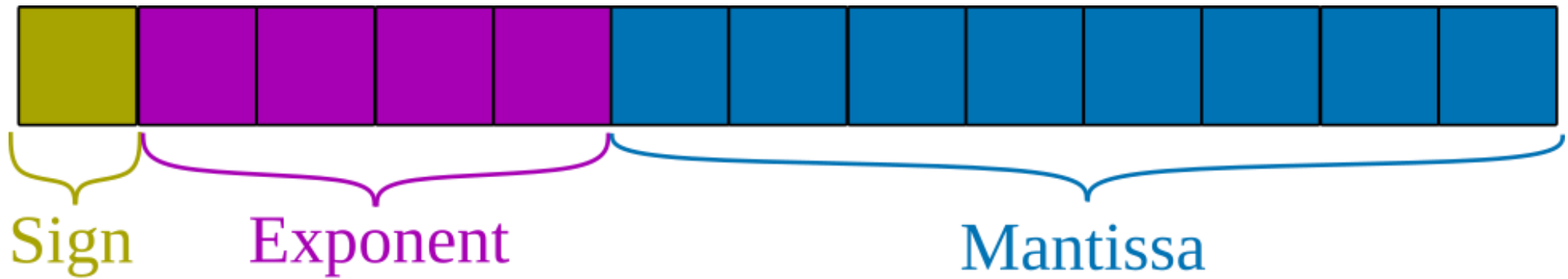
Outline :

I. Introduction

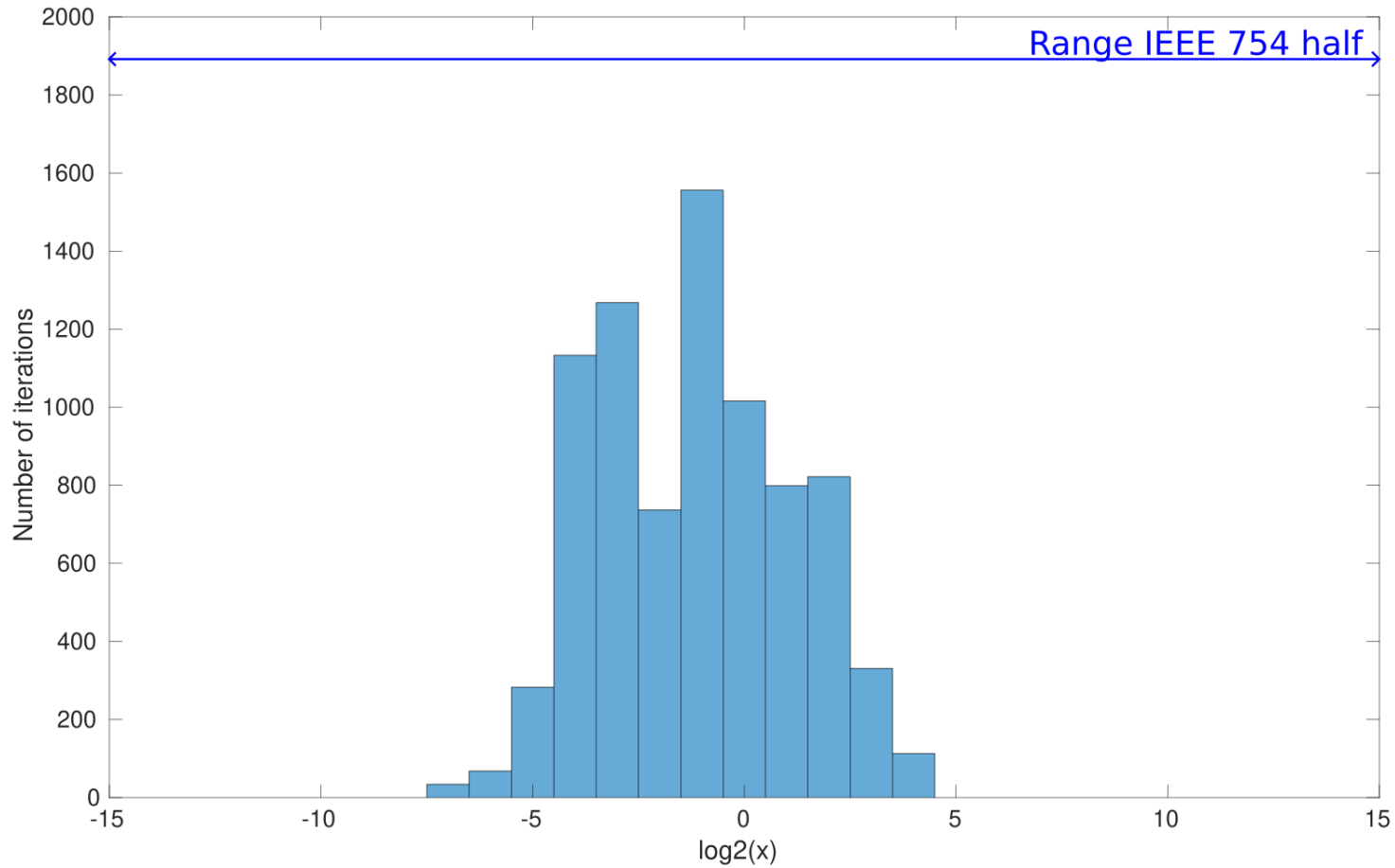
II. Optimization flow

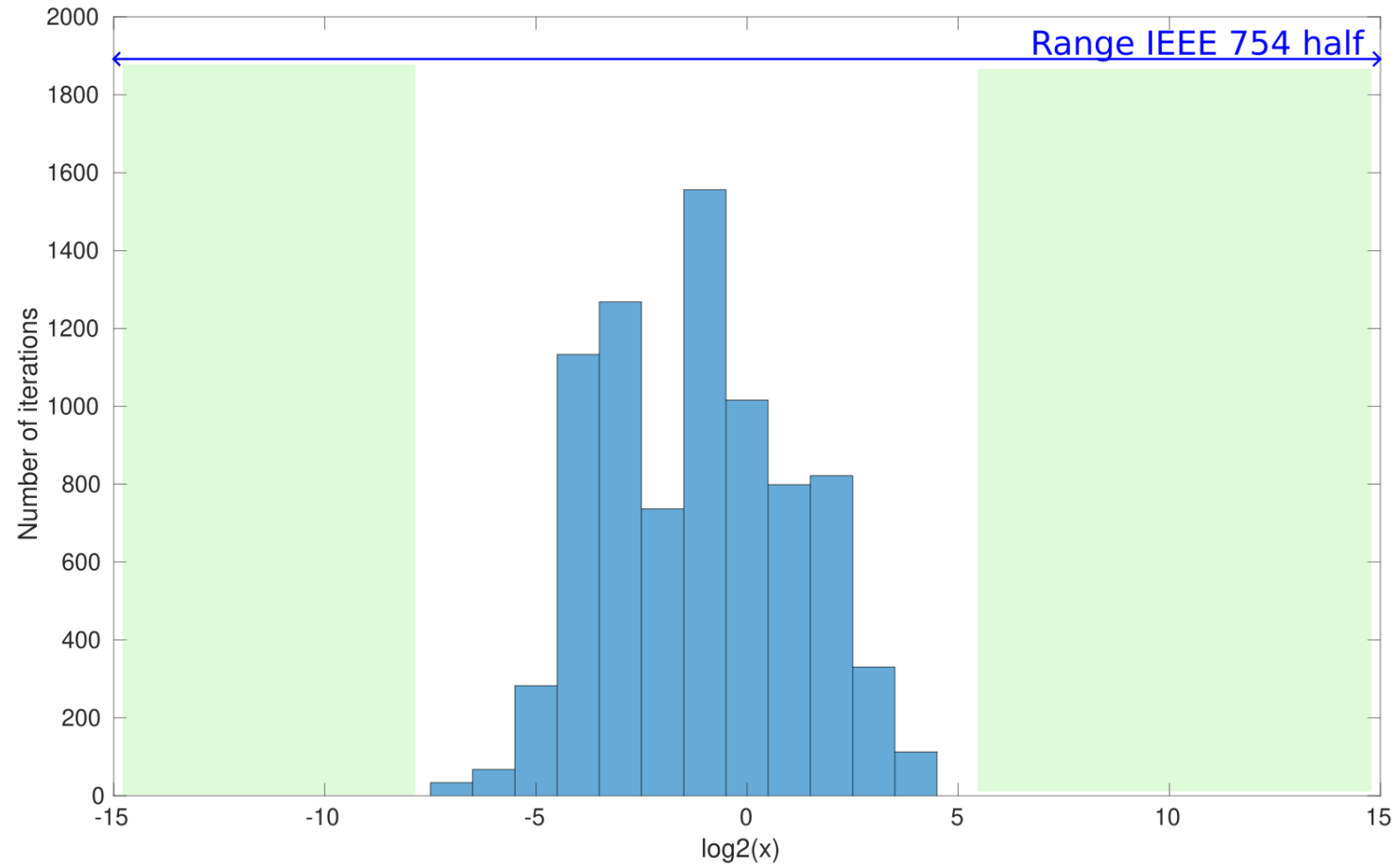
III. Proposed strategies

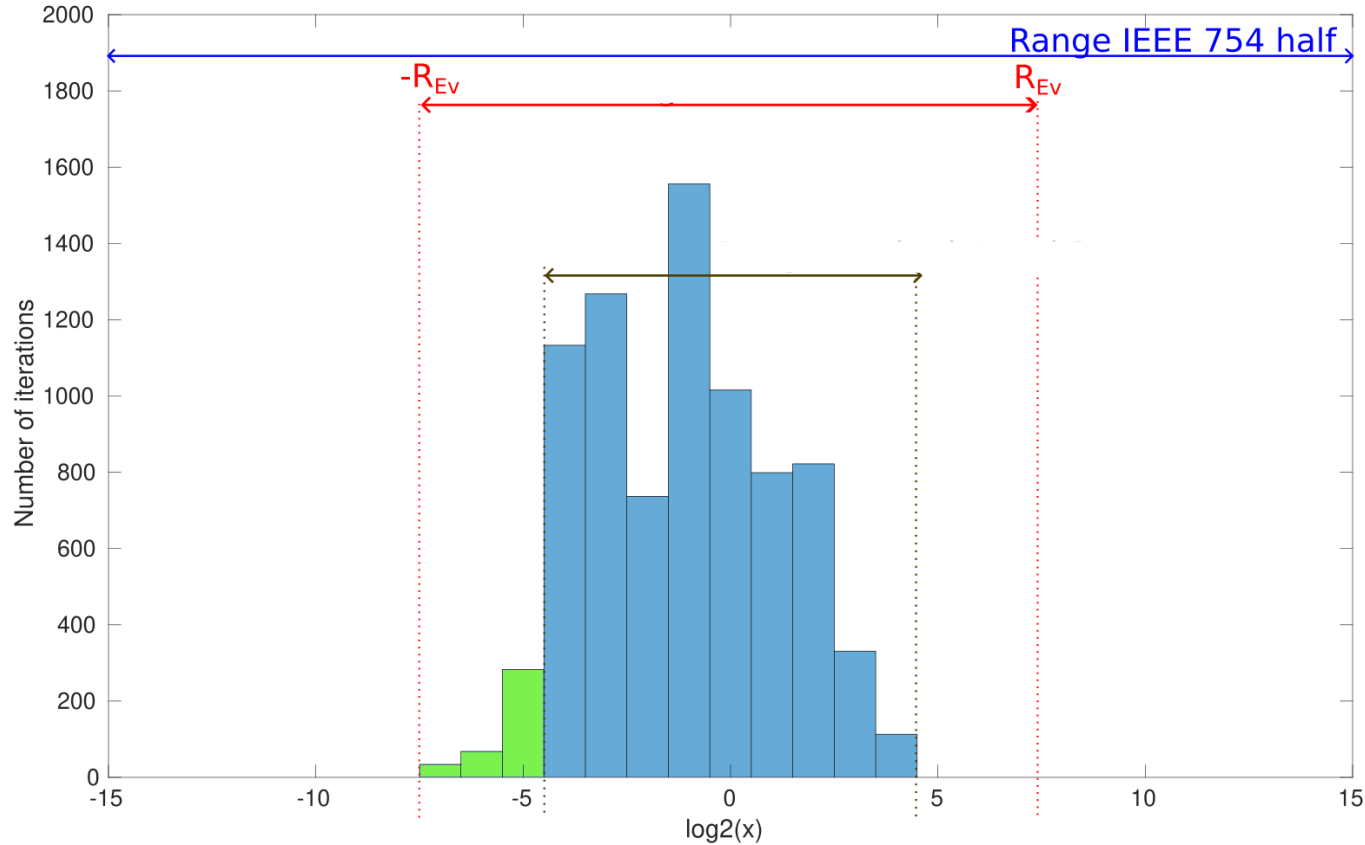
IV. Experiments



$$X = (-1)^S * (1 + M) * 2^{E-\Delta}$$

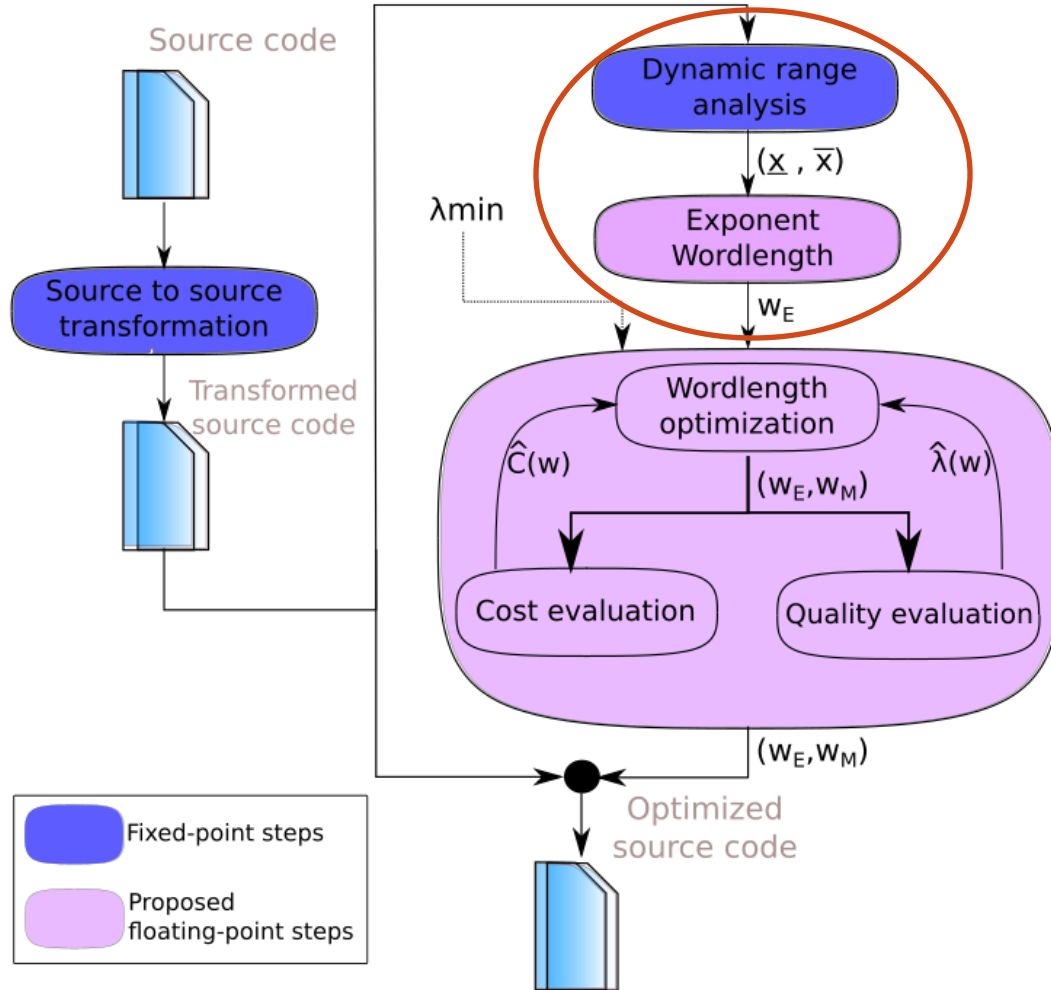






Number of exponents bits required :

$$W_{E_{Ri}} = \lceil \log_2 (R_{E_v}) \rceil + 1$$



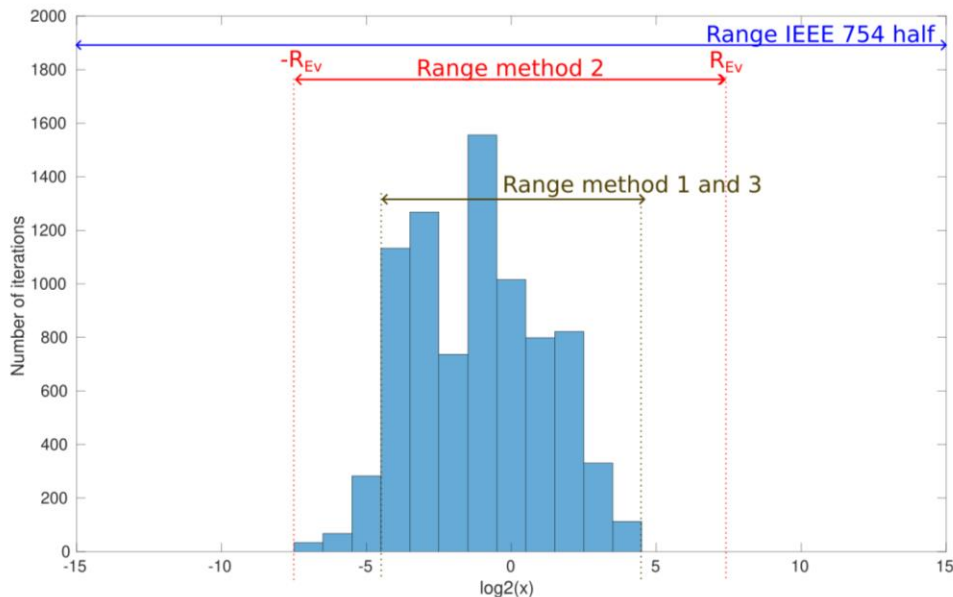
Reduced exponent wordlength

With :

- **Input** : $(W_{ER1}, W_{ER2}, \dots, W_{ERN}, W_{Mmax}, W_{Mmax}, \dots, W_{Mmax})$
- λ_{min} : The user minimal quality

Expected pros

- Explore a lot of useful configuration
- Explore all dimensions

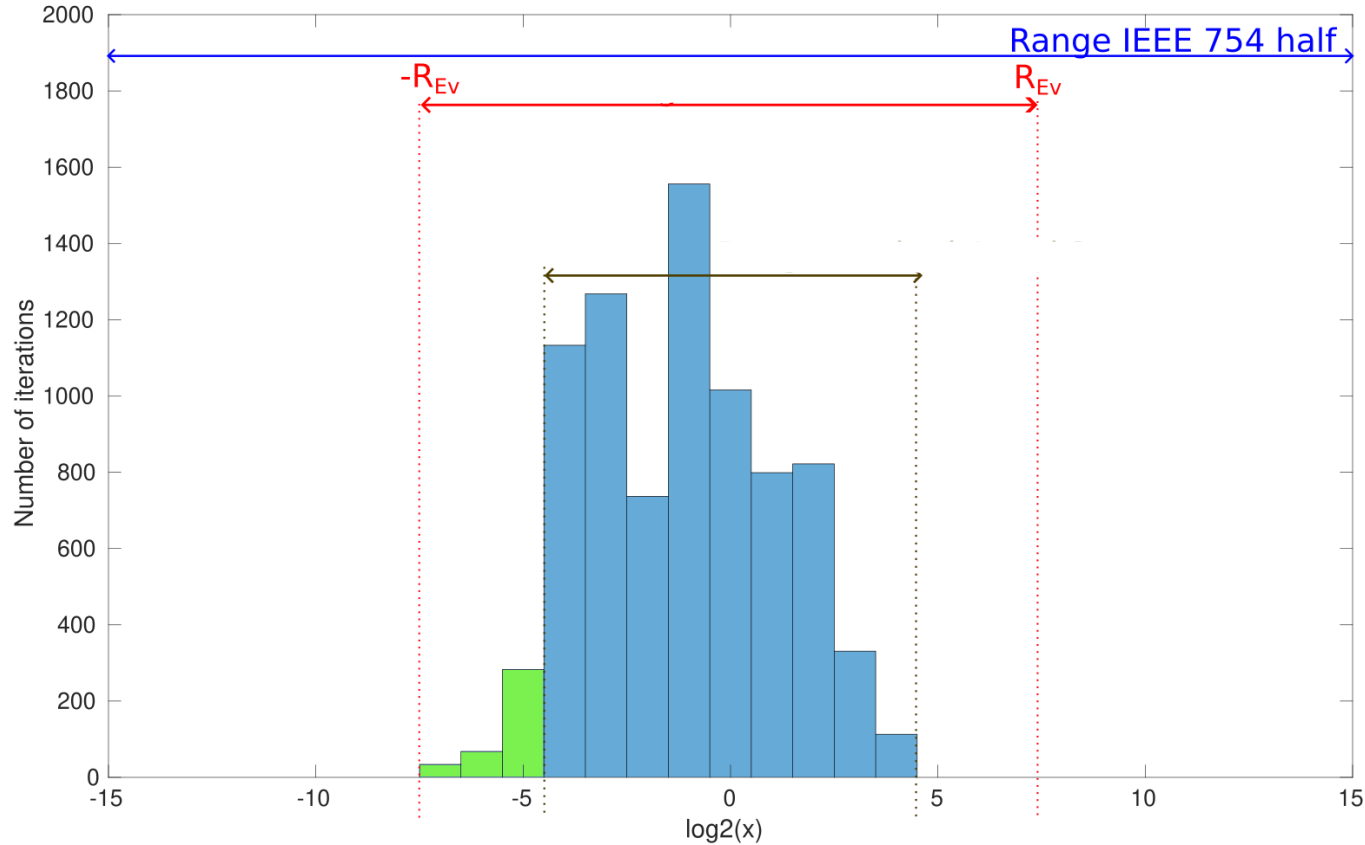


Expected cons

- Still a huge configuration number required :

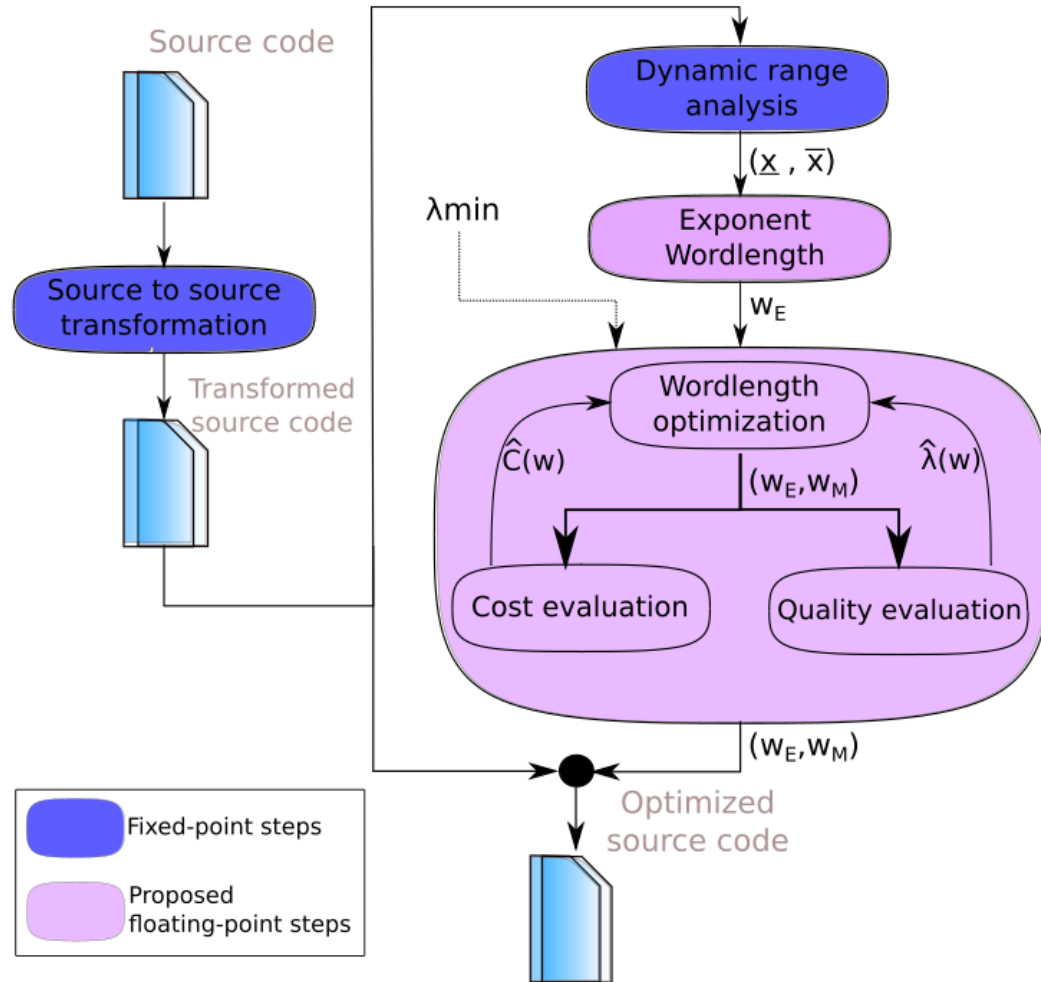
$$n_p = M^{N_v} * \prod_{i=0}^{N_v} w_{E_i}^{\max}$$

- Huge Optimization time



Number of exponents bits required :

$$W_{E_{Ri}} = \underbrace{\lceil \log_2 (R_{E_v}) \rceil} + 1$$



Only mantissa wordlength

With :

- **Input** : $(W_{M_{max}}, W_{M_{max}}, \dots, W_{M_{max}})$
- λ_{min} : The user minimal quality
- **All exponent size set to** $W_{E_{Ri}}$

Expected pros

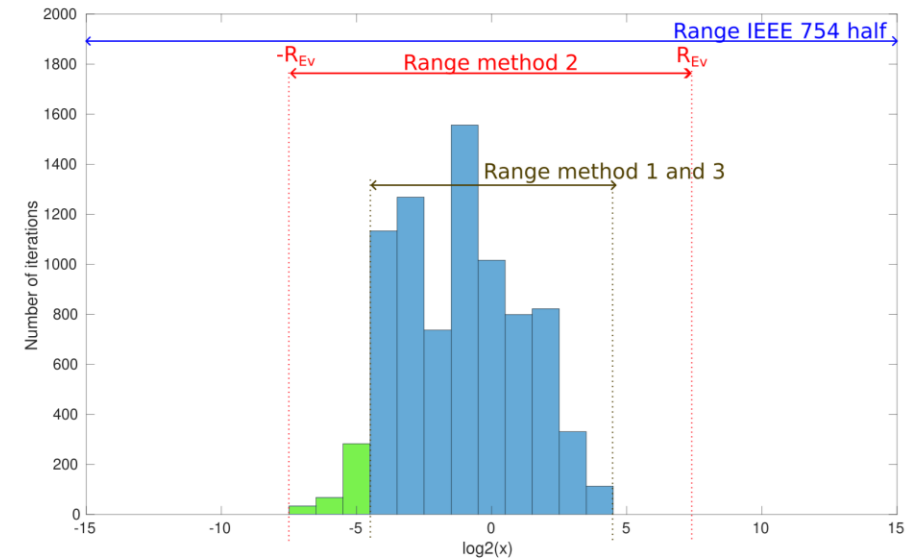
- Small configuration number required :

$$n_p = M^{N_v}$$

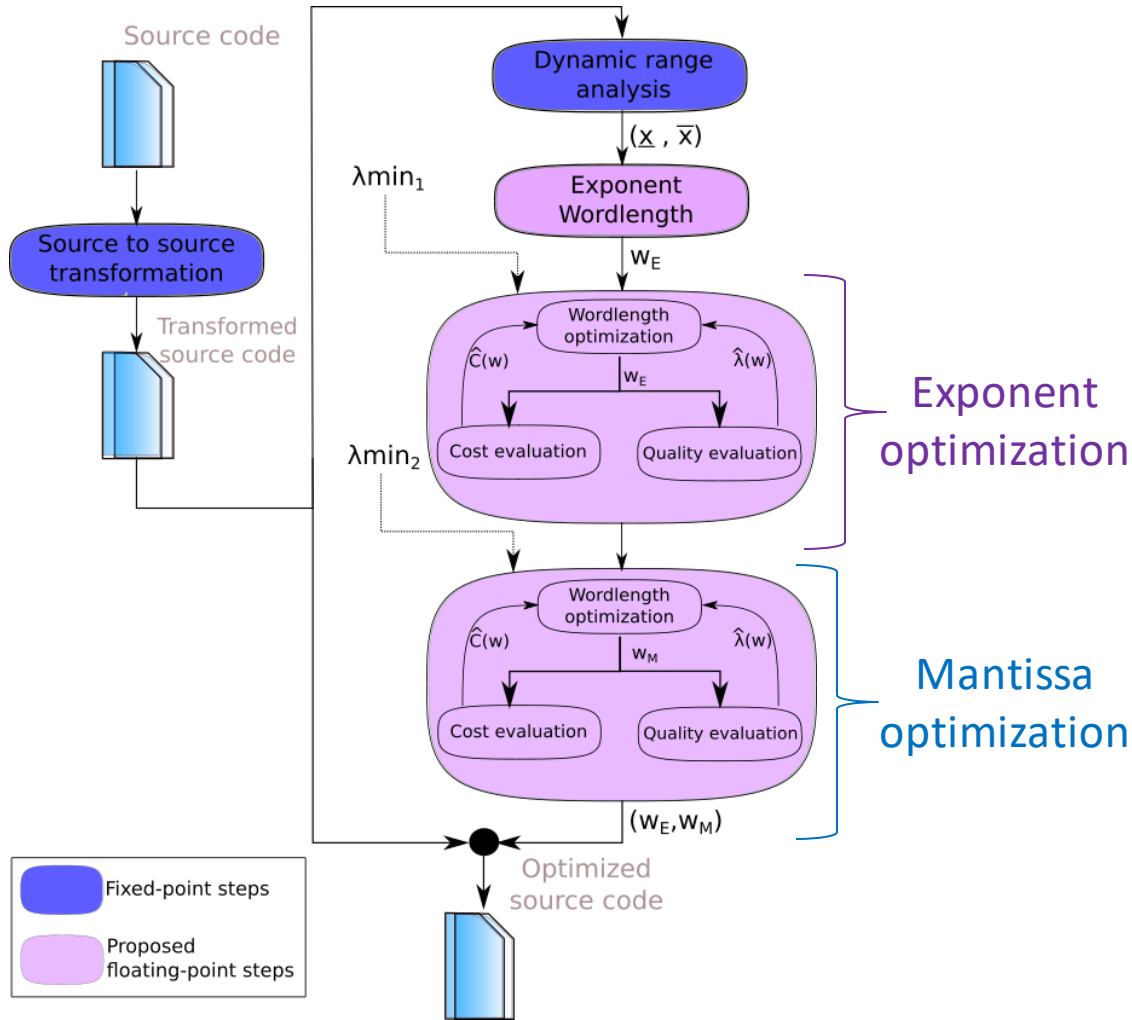
- Smaller Optimization time

Expected cons

- Doesn't try to reduce the exponent wordlength



- Lower quality of optimization



With :

- Input 1 : $(W_{E_{R1}}, W_{E_{R2}}, \dots, W_{E_{RN}})$
- λ_{min_1} : First user minimal quality
- All mantissa size set to $W_{M_{max}}$

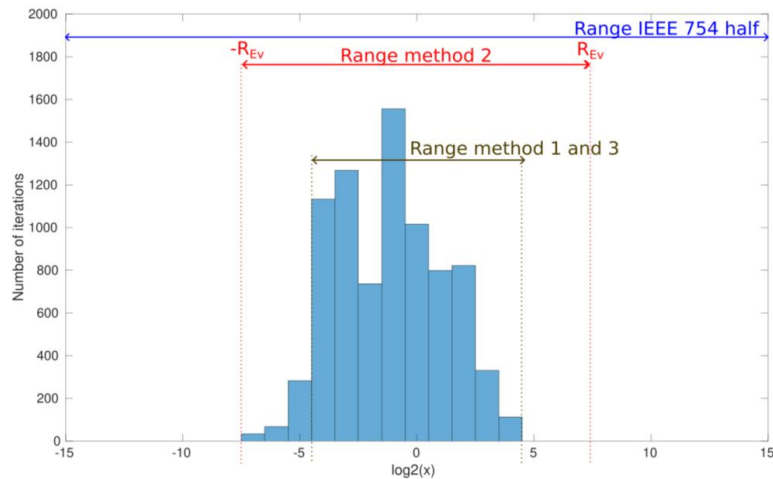
- Input 2 : $(W_{M_{max}}, W_{M_{max}}, \dots, W_{M_{max}})$
- λ_{min_2} : Second user minimal quality
- All exponent size set to W_{E_i} obtained with first optimization

Expected pros

- Middle ground for configuration number :

$$n_p = M^{N_v} + \prod_{i=0}^{N_v} w_{E_i}^{\max}$$

- Middle ground for the result quality



Expected cons

- Required and dependent of an hyperparameters : λ_{min_1}

Outline :

I. Introduction

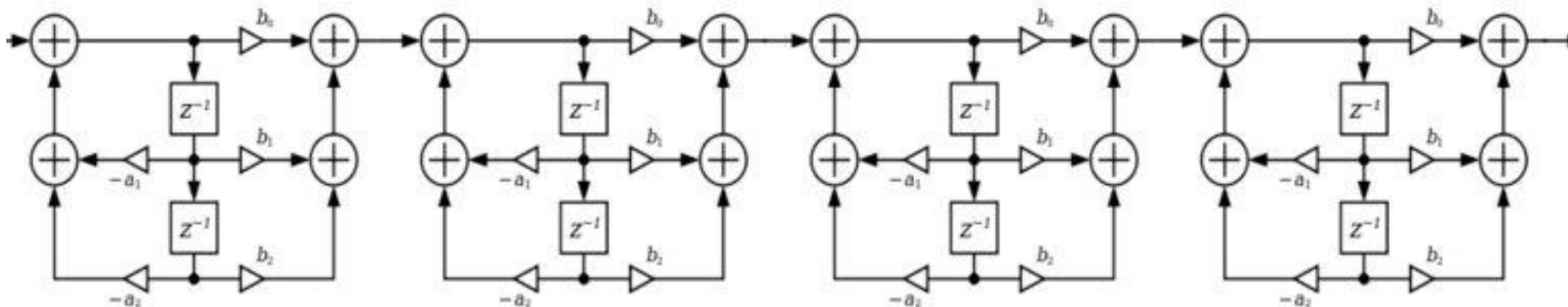
II. Optimization flow

III. Proposed strategies

IV. Experiments

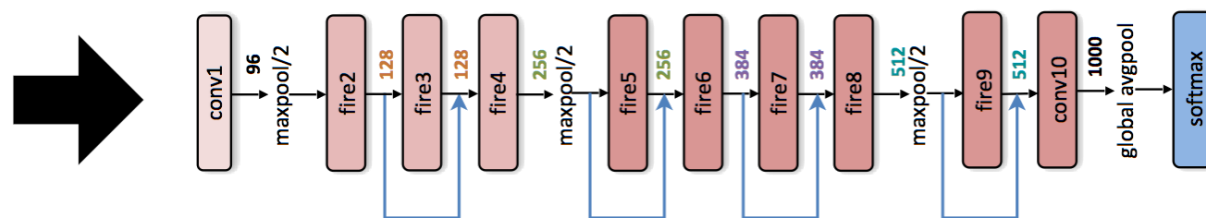
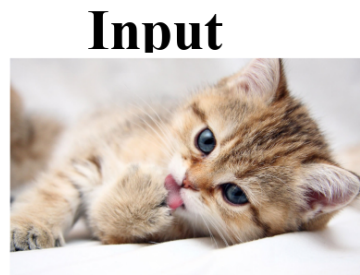
Experiment 1: Infinite Impulse Response filter

- Quality metric : Signal to Quantization Noise Ratio (SQNR)
- Number of variable to optimize : 7

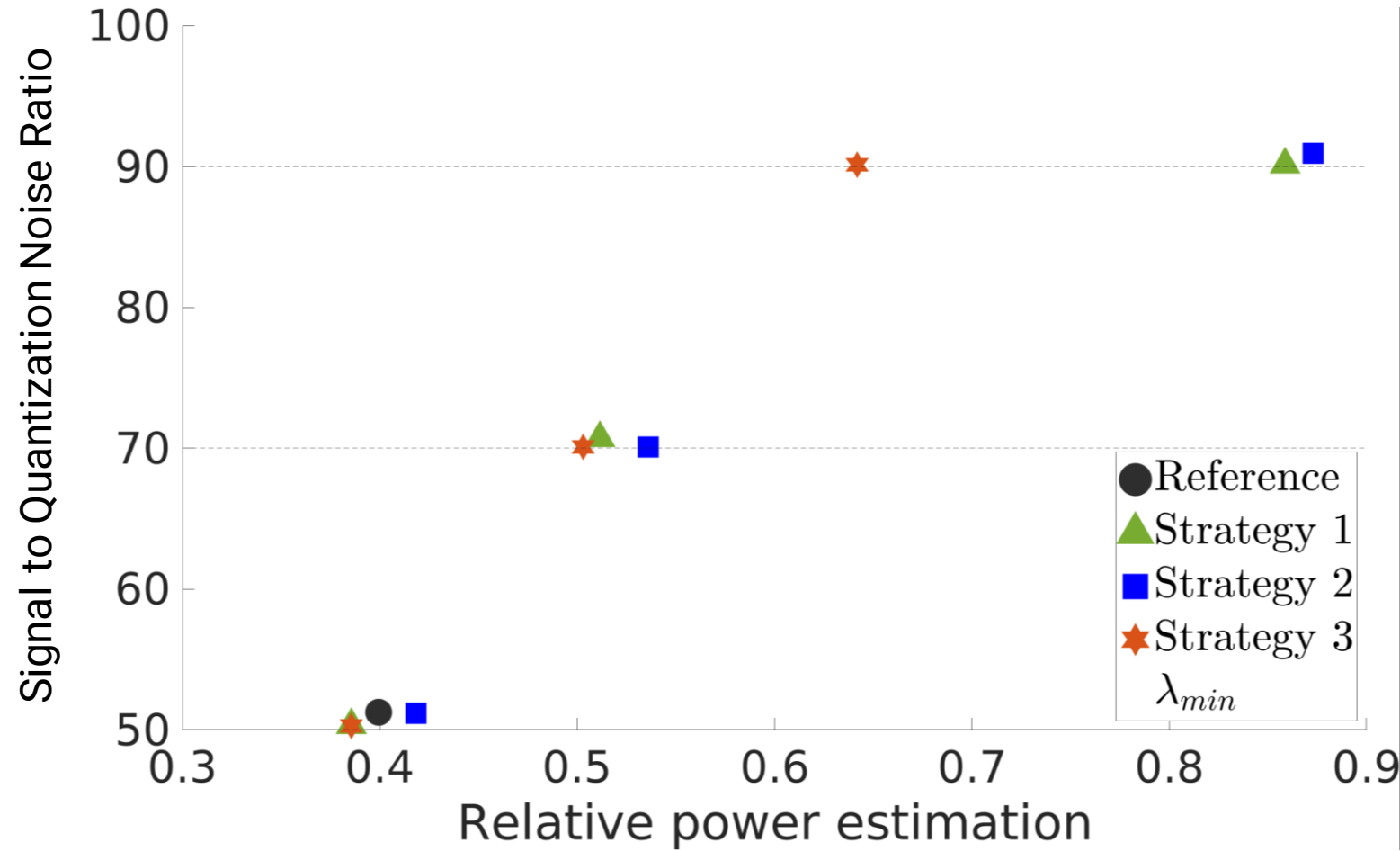


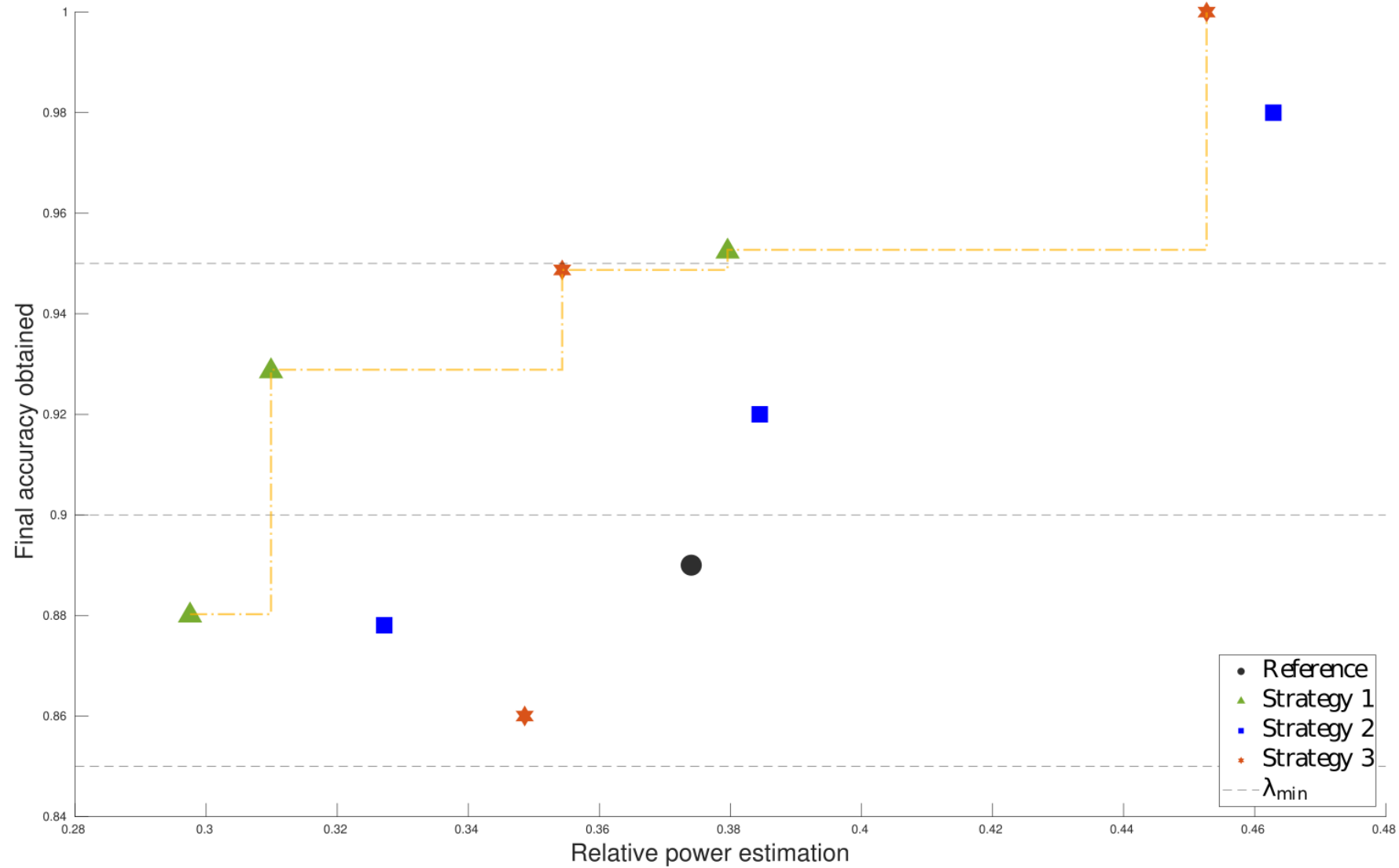
Experiment 2: Squeezenet

- Quality metric : proportion of similar top 1 classes detected
- Number of variable to optimize : 15

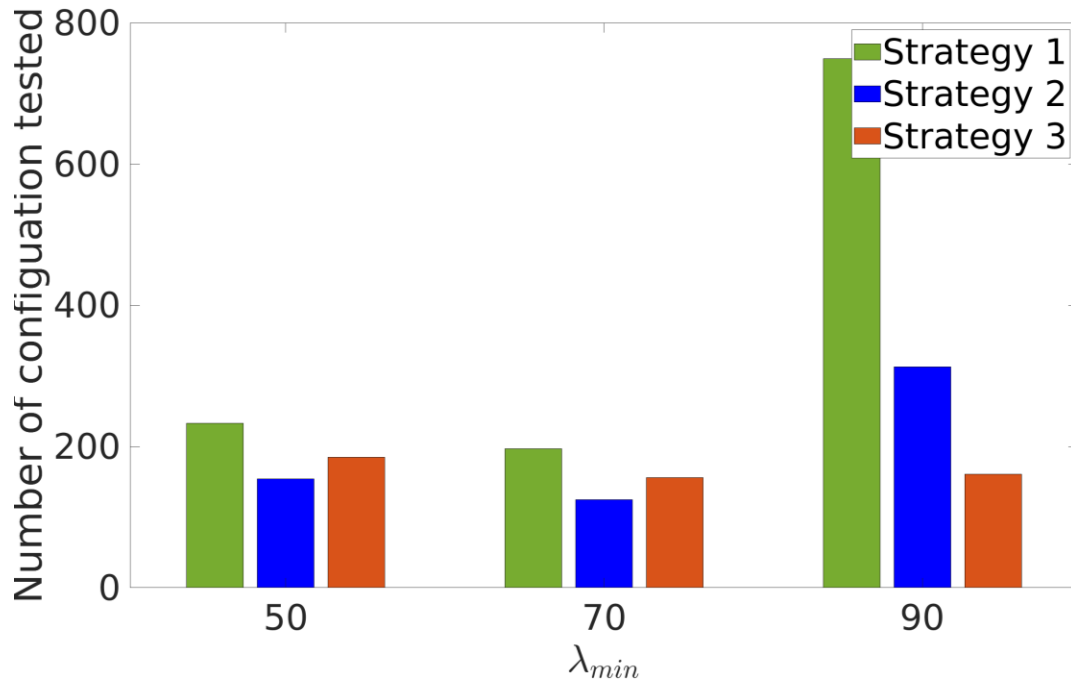


1.cute cat: 91%
2.dog : 52%
3.fly : 12%

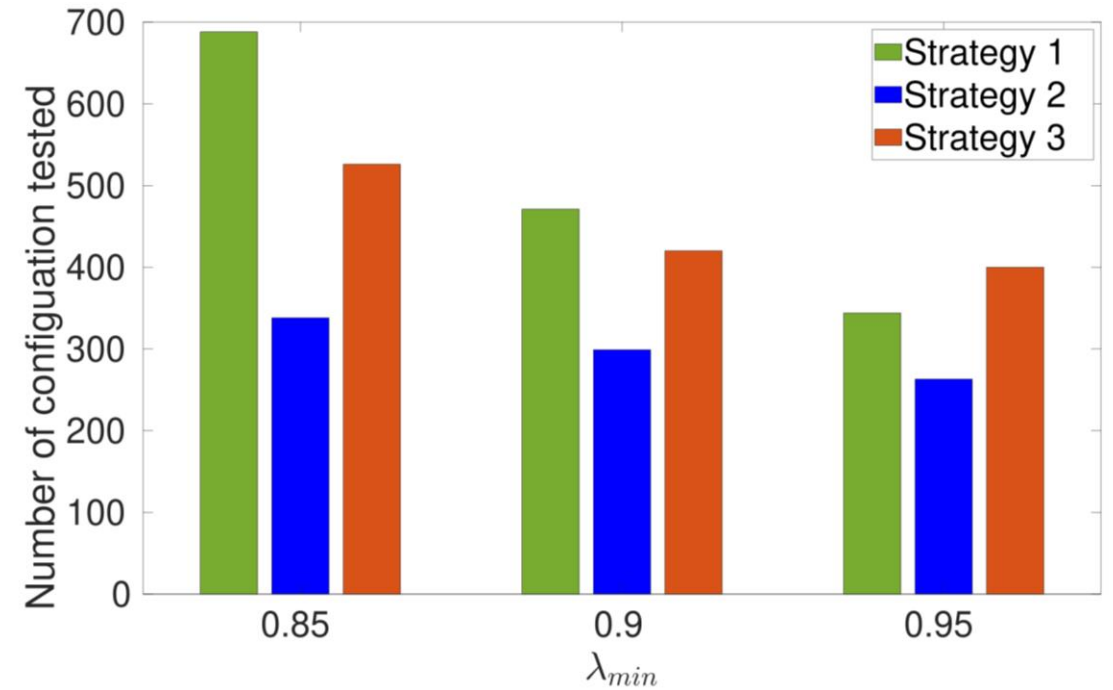




Infinite Impulse Response filter



Squeezenet



- Proposed :
 - New custom floating-point refinement
 - Three strategies to optimize exponent and mantissa wordlength
 - Allow to follow a user-defined quality constraint
- Find improvement in terms of power and memory consumption compared to half floating-point
- Perspective :
 - Increase the quality of optimization
 - Implement memory cost reading/storing energy in power estimation