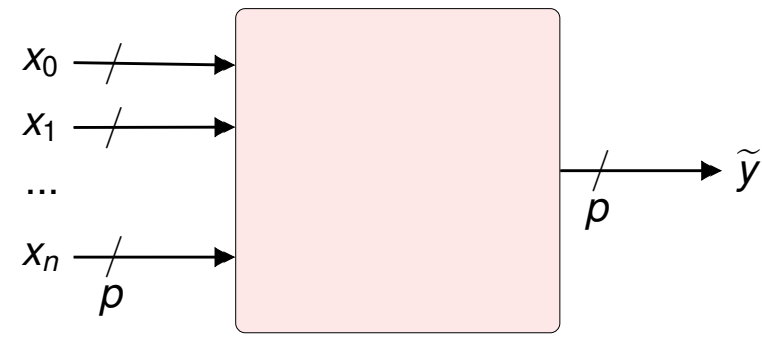


Precision vs. Accuracy in Hardware – The FloPoCo Philosophy

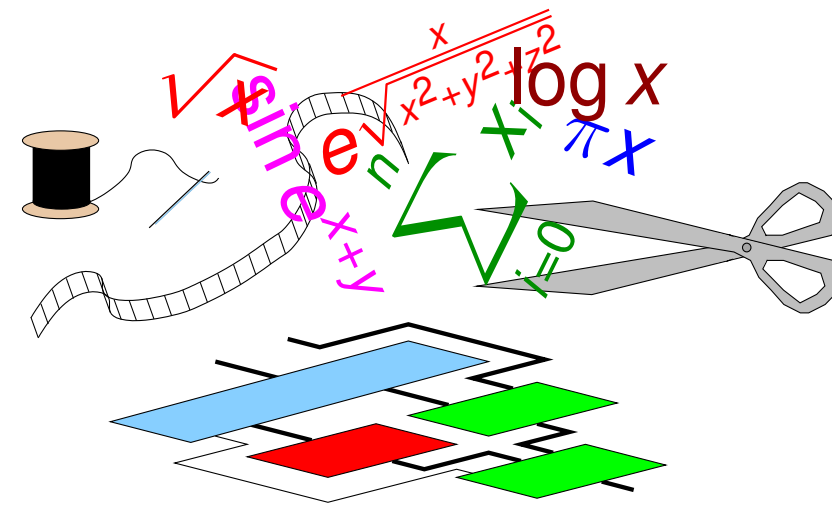
Computing Just Right



Why compute the result with **more accuracy** than the output can hold?
⇒ wasted accuracy that nobody can see

Why compute with **less accuracy** than your output can hold?
⇒ part of the result will be just noise

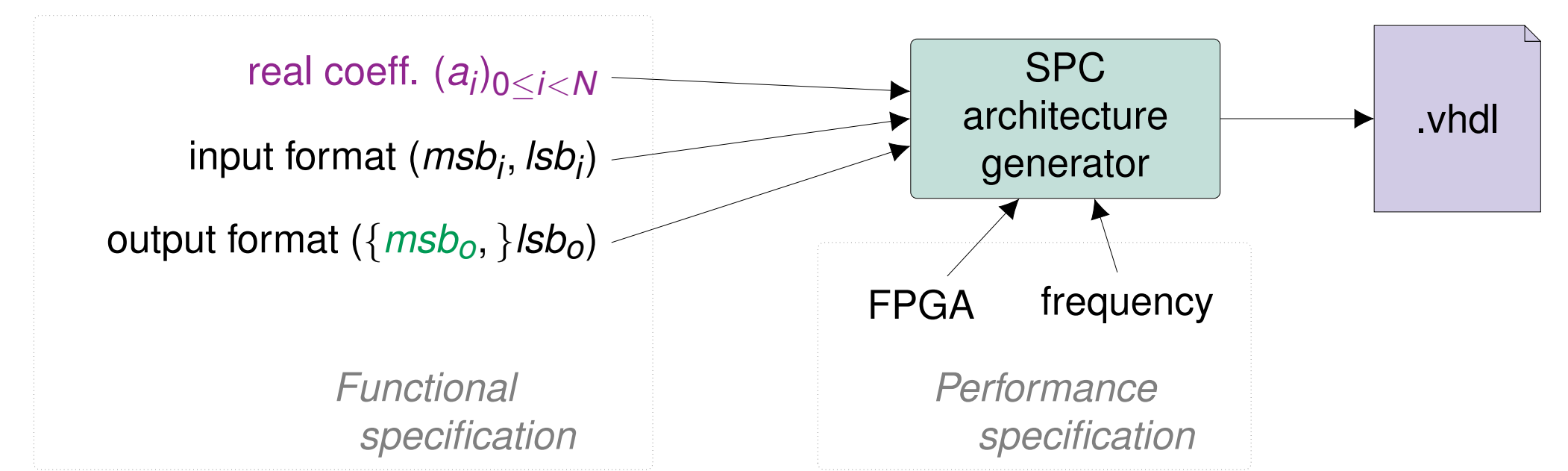
The FloPoCo Arithmetic Core Generator



Accurate results with **minimal resources**:

- generator and library for custom FPGA targets
- explore the flexibility of an FPGA

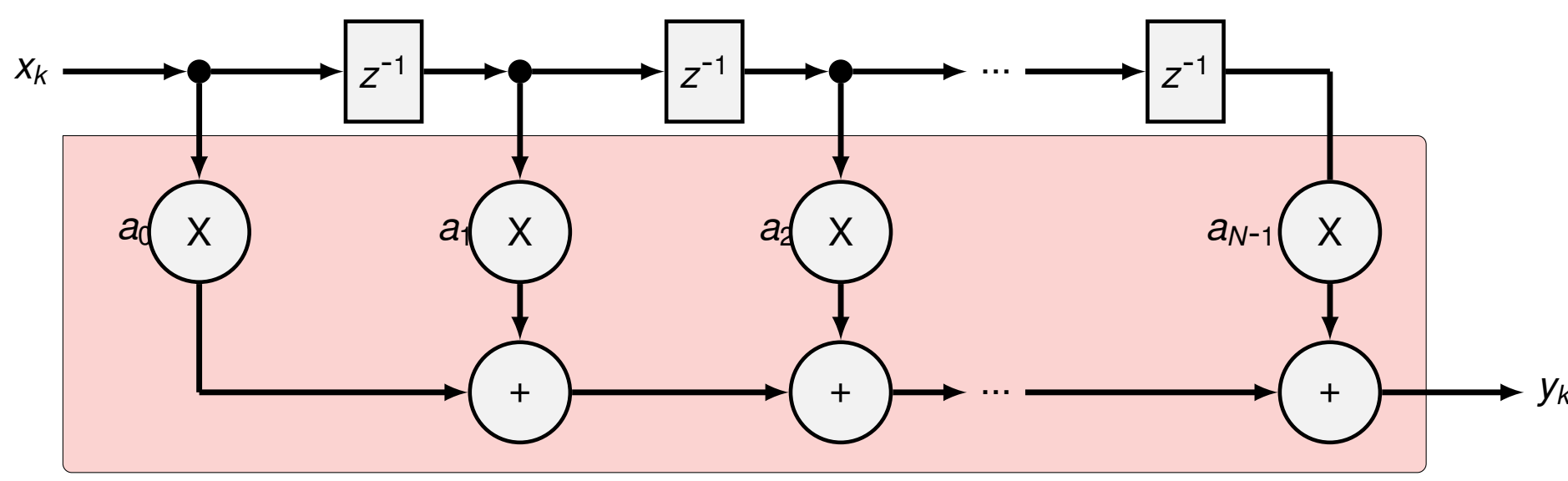
Example Operator – An Automatic Filter Generator



Output accuracy defined by the output format.

Implementing FIR Filters in Hardware

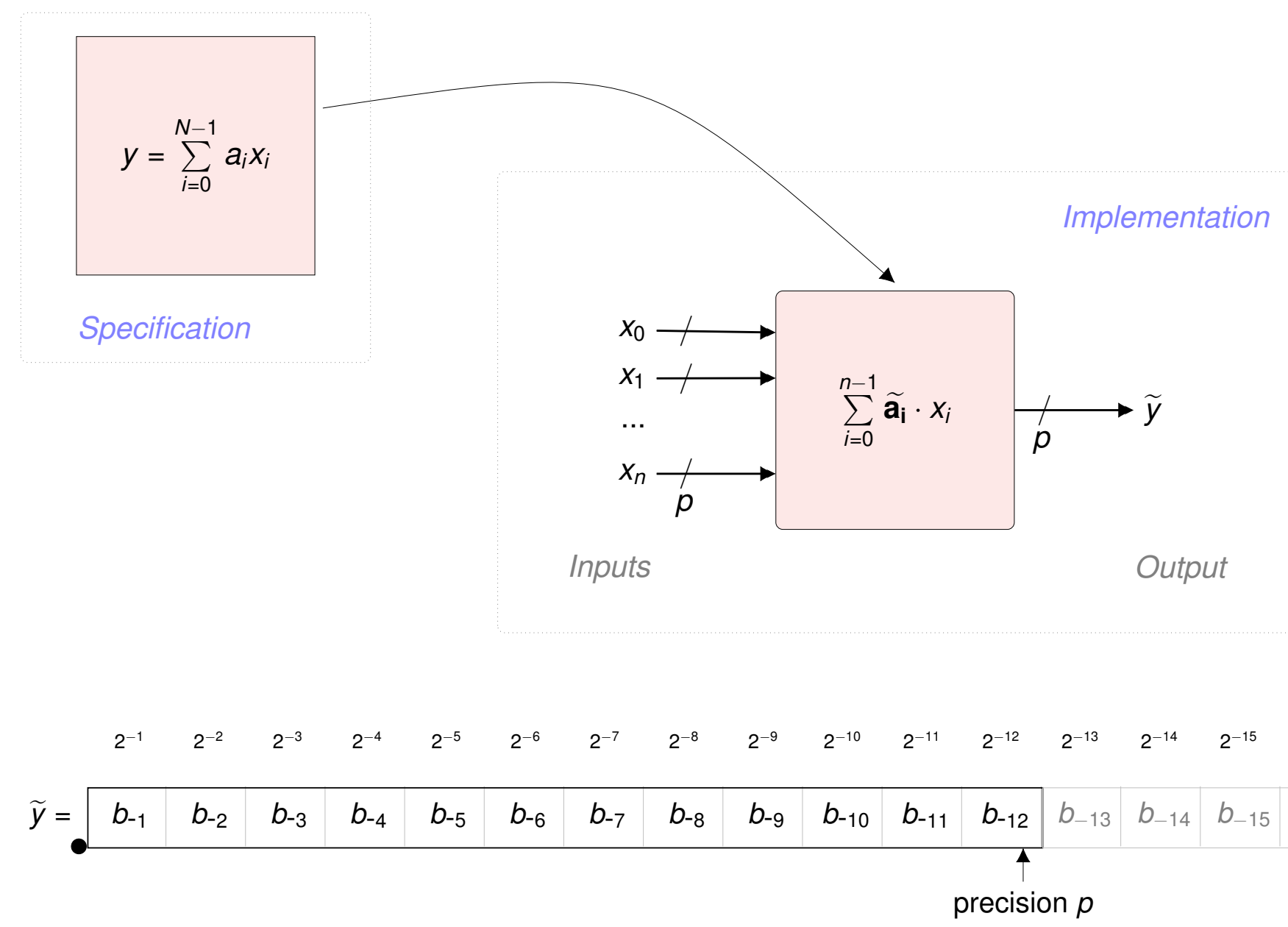
A classical solution



Each element y_k is a **sum of products**.

$$y_k = \sum_{i=0}^{N-1} a_i x_{k-i}$$

Sums of Products: From Theory to Implementation



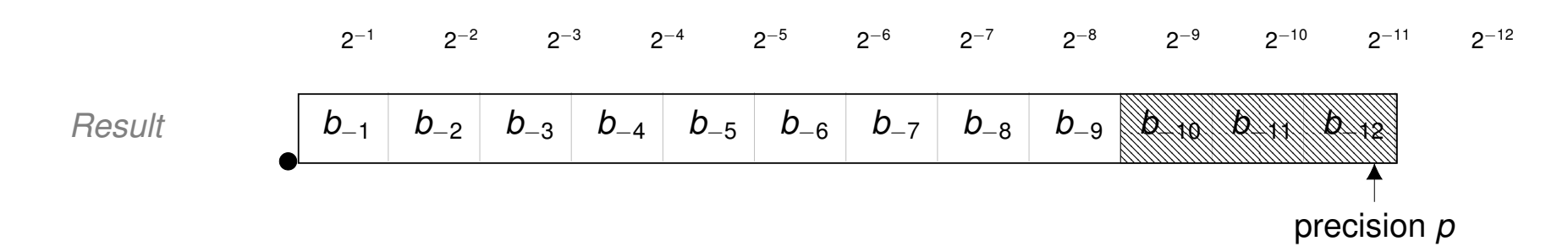
Computing with precision p ?

What we **intended** to compute: the mathematical result

$$y = \sum_{i=0}^{N-1} a_i x_i$$

What we **actually** computed:

$$\tilde{y} = \sum_{i=0}^{N-1} \text{round}_p(\tilde{a}_i x_i)$$



The error:

$$\epsilon = \tilde{y} - y$$

Computing Just Right A Sum of Products by Constants

Determine the Sum's Most Significant Bit (MSB)

The MSB of $a_i x_i$

- x_i bounded (fixed-point number)
- a_i known

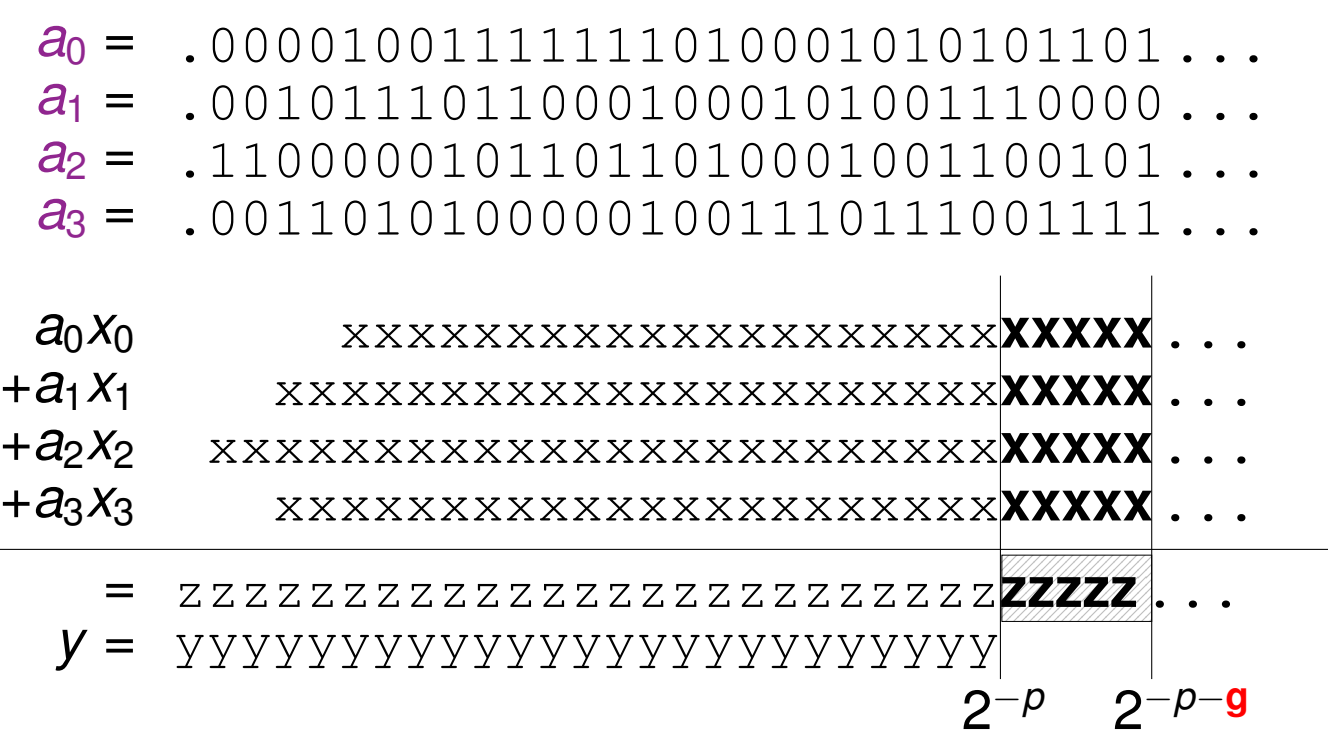
$$msb_{a_i x_i} = \lceil \log_2(|a_i| \text{val}_{max}(x_i)) \rceil$$

The MSB of the sum

- $a_i x_i$ bounded

$$msb_0 = msb_y = \lceil \log_2(\sum_{i=0}^{N-1} |a_i| \text{val}_{max}(x_i)) \rceil$$

Determine the Sum's Least Significant Bit (LSB)



Suppose we use perfect multipliers: $\epsilon_{mult} < 2^{-p-g-1}$

- after the sum:

$$\epsilon_{total} = \sum_{i=0}^N \epsilon_{mult} + \epsilon_{final_rounding} < N \cdot 2^{-p-g-1} + 2^{-p-1}$$

Need for larger intermediary precision

- **g guard bits**
- such that errors accumulate in the guard bits
⇒ $g = \lceil \log_2(N) \rceil$

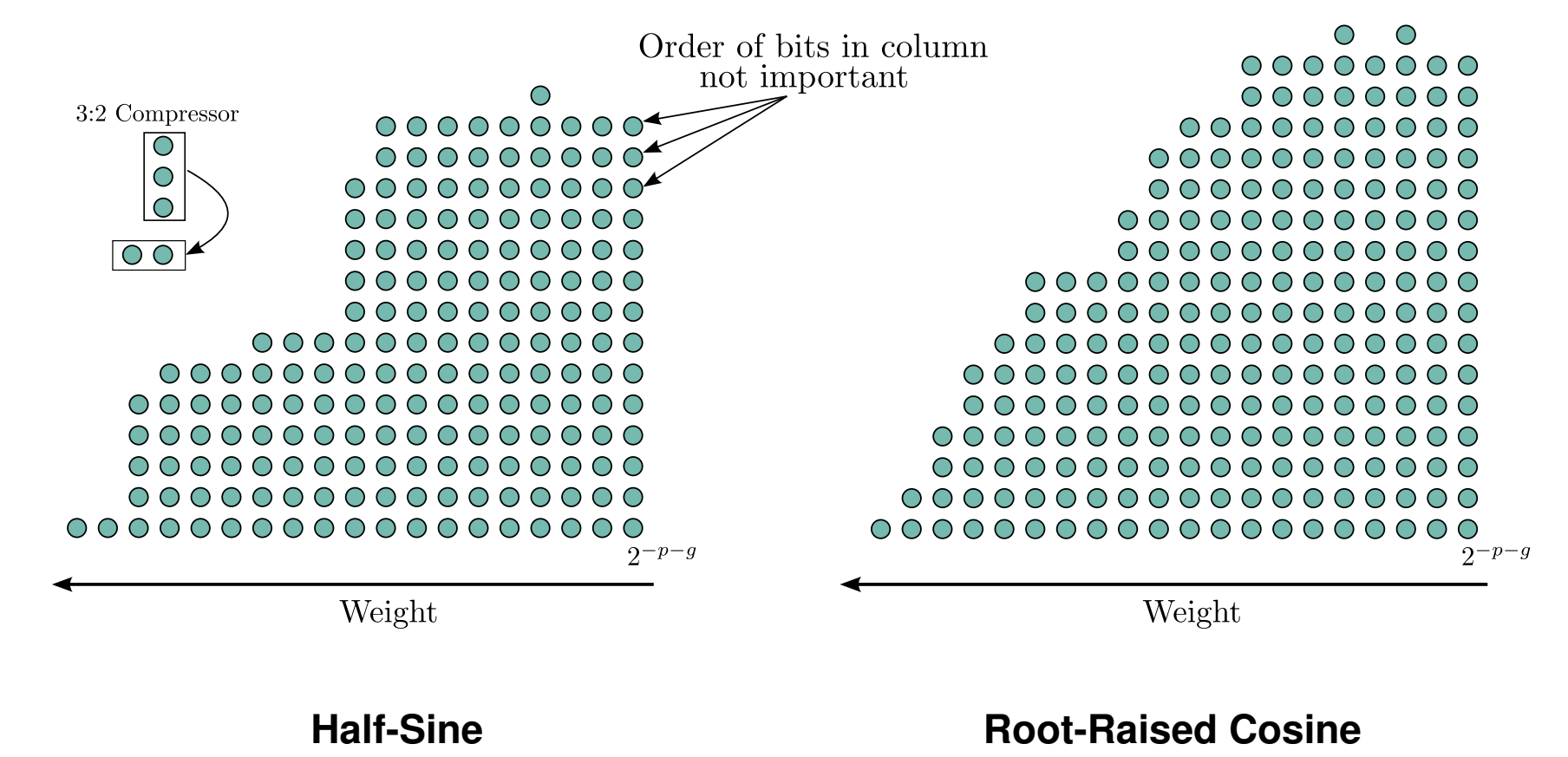
Computing the Sum

- **Bit-heaps** in FloPoCo

- bit-parallel computation of the sum
- generalization of bit arrays used in multipliers

- Example:

8-tap, 12-bit FIR filters used in the Zigbee protocol

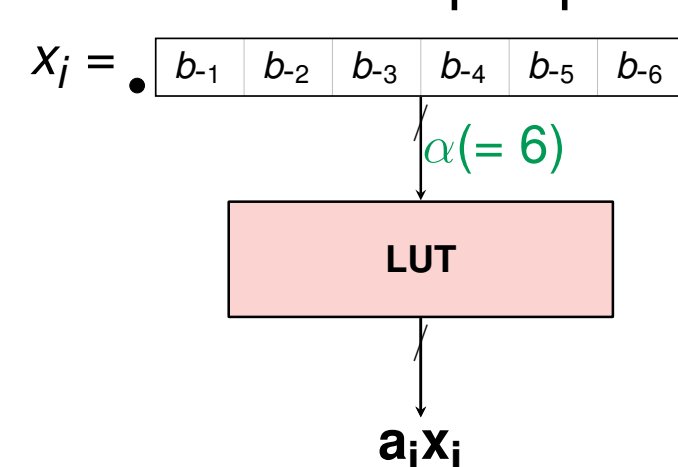


Multiplying by a Constant in a FPGA

Perfect Constant Multipliers

Basic FPGA block: look-up table (LUT) of α input bits

- tabulate all the values of $a_i x_i$
- ... **correctly rounded** to the output precision



Perfect fit for small sizes: x_i on α bits ⇒ 1 LUT/output bit

Doesn't scale:

- 2^k LUT/output bit for x_i on $(\alpha + k)$ bits

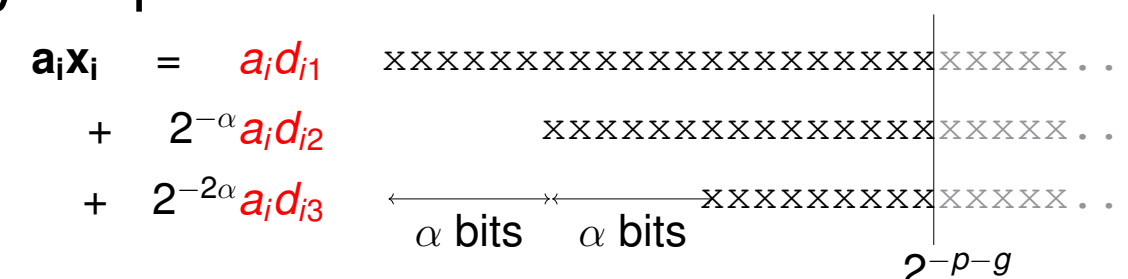
KCM Multipliers

When $x_i > \alpha$ bits: cut it into chunks d_{ik} of α bits

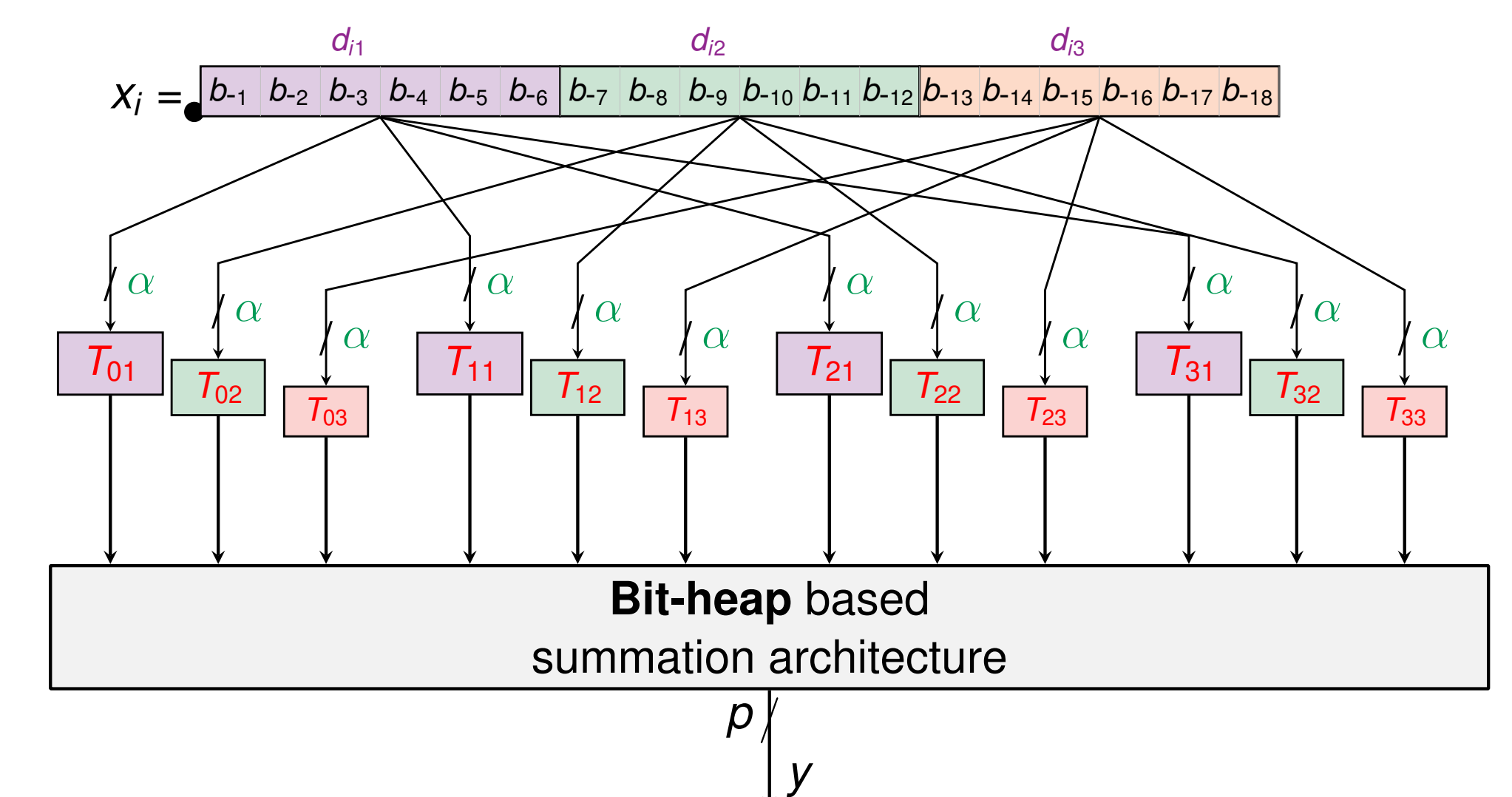
$$x_i = \sum_{k=1}^n 2^{-k\alpha} d_{ik} \quad \text{where } d_{ik} \in \{0, \dots, 2^\alpha - 1\}$$

$$\Rightarrow a_i x_i = \sum_{k=1}^n 2^{-k\alpha} a_i d_{ik}$$

Each $a_i d_{ik}$ tabulated, 1 LUT/output bit
How many output bits?



Sum of Products Using KCM Multipliers



Example – Root-Raised Cosine Filter

Method	p	Latency	LUTs	Error Bound	Output Bits
Naïve method	$p = 12$	5.9 ns (170 Mhz)	444 LUT	$\bar{\epsilon} > 2^{-9}$	$y_1, y_0, y_{-1}, y_{-2}, y_{-3}, y_{-4}, y_{-5}, y_{-6}, y_{-7}, y_{-8}, y_{-9}, y_{-10}, y_{-11}, y_{-12}$
Proposed method	$p = 12$	4.4 ns (227 Mhz)	564 LUT	$\bar{\epsilon} < 2^{-12}$	$y_1, y_0, y_{-1}, y_{-2}, y_{-3}, y_{-4}, y_{-5}, y_{-6}, y_{-7}, y_{-8}, y_{-9}, y_{-10}, y_{-11}, y_{-12}$
Proposed method	$p = 9$	4.12 ns (243 Mhz)	380 LUT	$\bar{\epsilon} < 2^{-9}$	$y_1, y_0, y_{-1}, y_{-2}, y_{-3}, y_{-4}, y_{-5}, y_{-6}, y_{-7}, y_{-8}, y_{-9}$