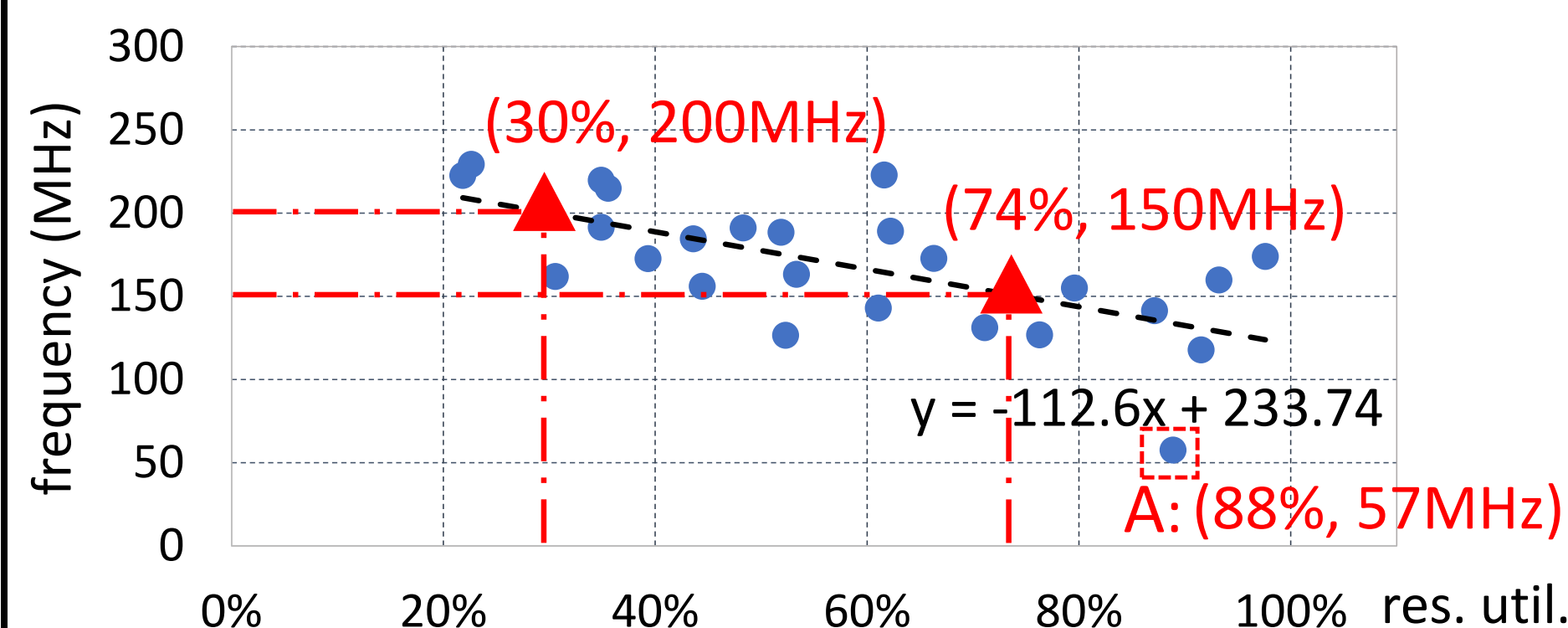


Motivation

◆ Motivation

- Frequency decreases as design size scales out in HLS accelerators
- Severe frequency degradation:
 - 30% chip resource -> 200MHz
 - 74% chip resource -> 150MHz
 - 90% chip resource -> 132MHz
 - Extreme case: FFT, 88% area -> 57MHz



Achilles' heel

- Identify four common collective communication and computation pattern: scatter, gather, broadcast and reduce.

- One-to-all or all-to-one **on-chip data movement**
- Wirelength scales up -> critical path delay
- Patterns are used in most, if not all, accelerators

Benchmark	Domain	Scatter	Gather	Broadcast	Reduce
AES	Encryption	✓*	✓*	✓	
FFT	Signal	✓	✓*		
GEMM	Algebra	✓	✓*	✓*	
KMP	String	✓		✓	✓*
NW	Bioinfo.	✓	✓		
SPMV	Algebra	✓	✓*	✓	
STENCIL	Image	✓	✓*	✓	
VITERBI	DP	✓	✓		

Checkmark ✓ represents the design has the pattern.
A star * represents that a critical path lies in the pattern.
For broadcast, GEMM uses bc_in_compute while others use bc_by_copy.

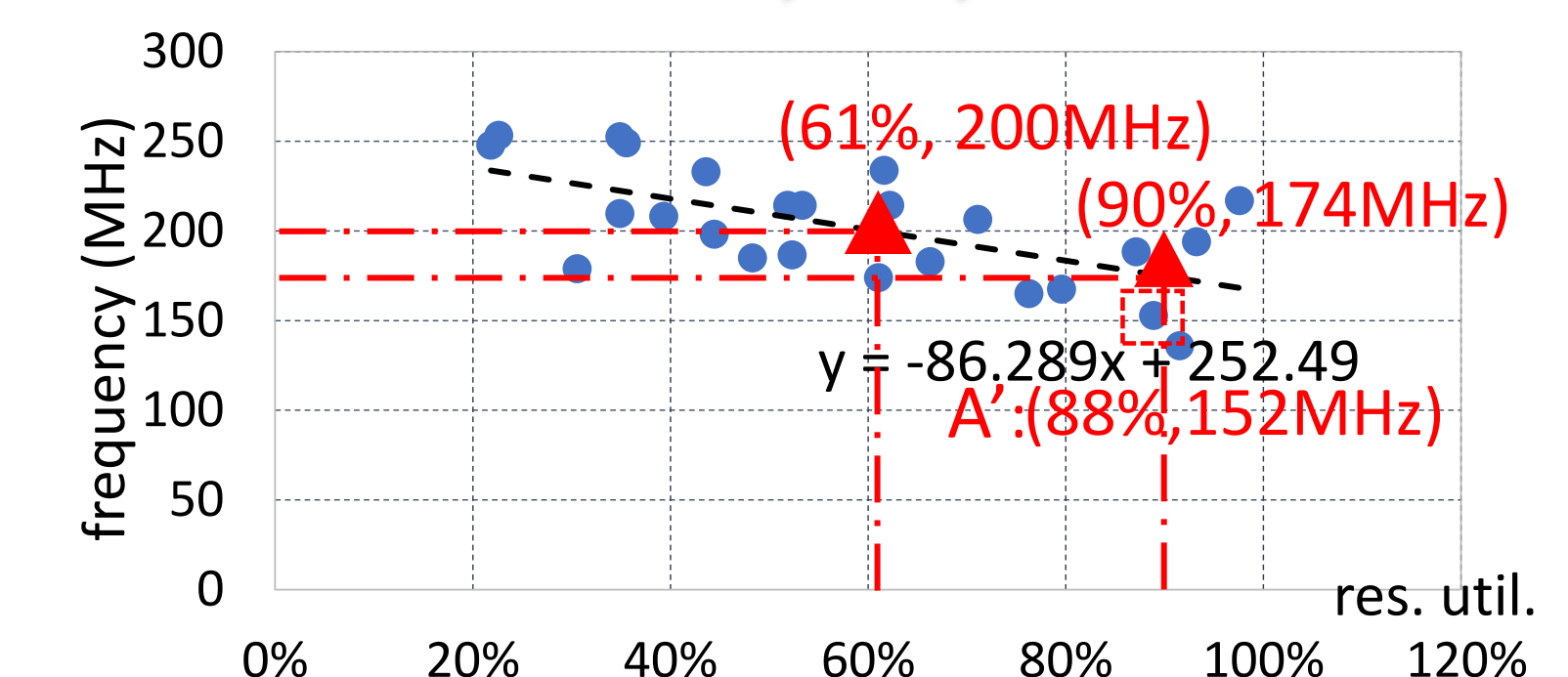


Ouch! Painful -> Latte

“Caffè” Latte

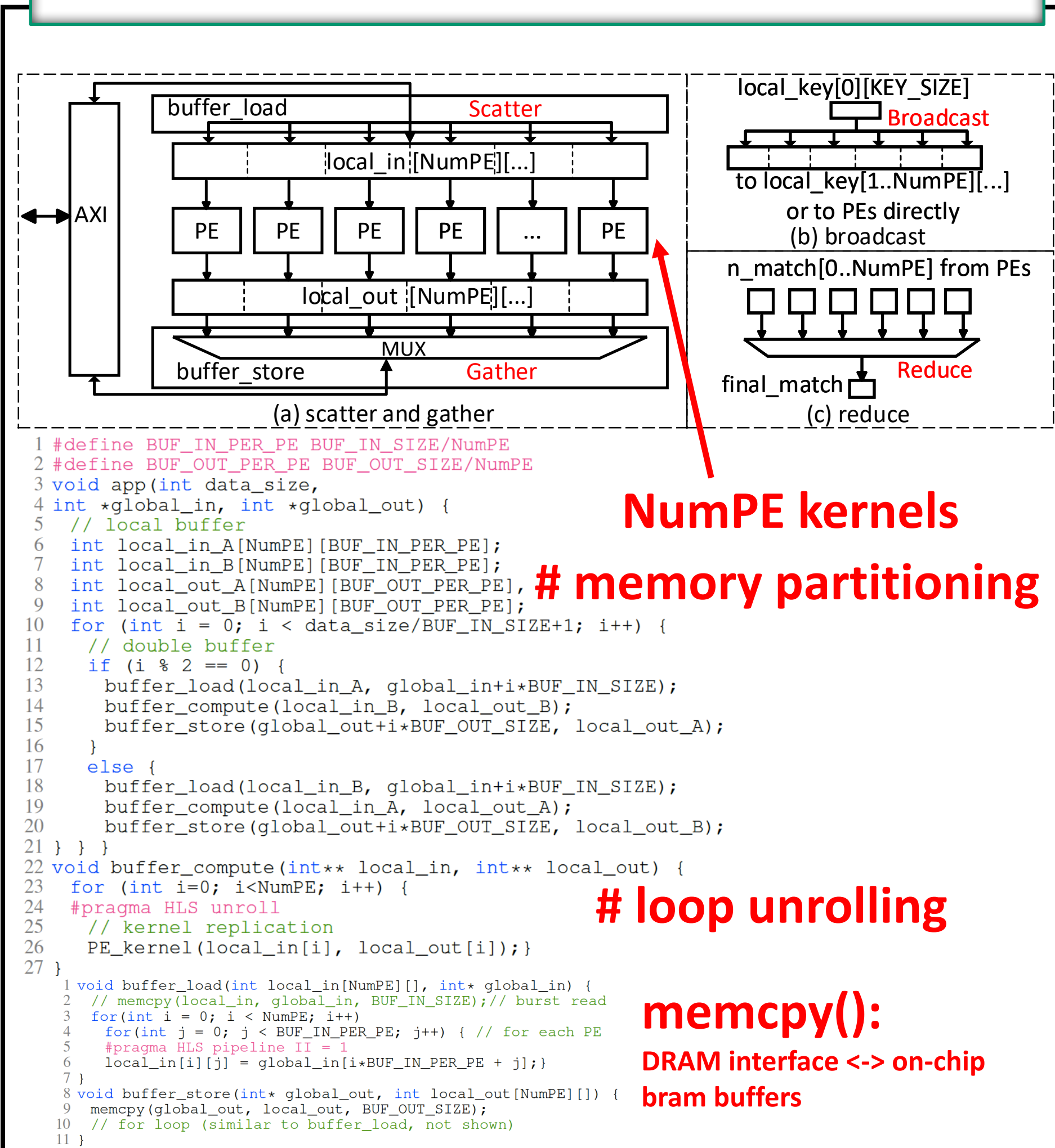
◆ Latte boosts frequency

- Latte-optimized design improves timing of the baseline by 1.50x with only 3.2% LUT overhead
- 30% chip resource, 200MHz -> 227MHz
- 74% chip resource, 150MHz -> 189MHz
- 90% chip resource, 132MHz -> 175MHz
- Extreme case: FFT, 88%, 57MHz -> 152MHz



How does Latte work?

Common Practice Accelerator



NumPE kernels

memory partitioning

loop unrolling

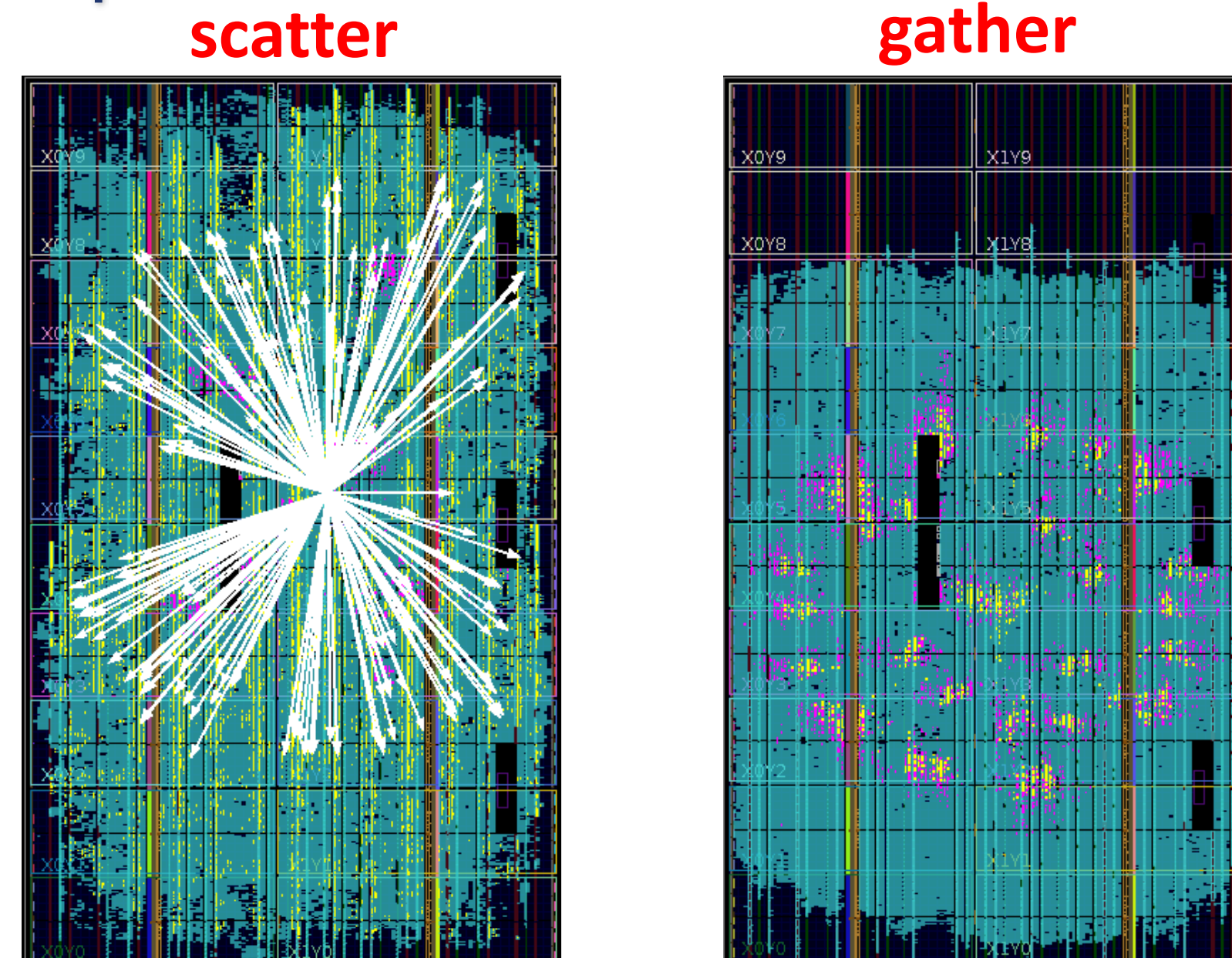
memcpy():

DRAM interface <-> on-chip
bram buffers

Chip Layout Look Like?

◆ Scattered distribution of local buffers (PEs)

- HLS optimistically estimate memcpy() function wire delay **without considering the locality of partitioned local buffer banks**

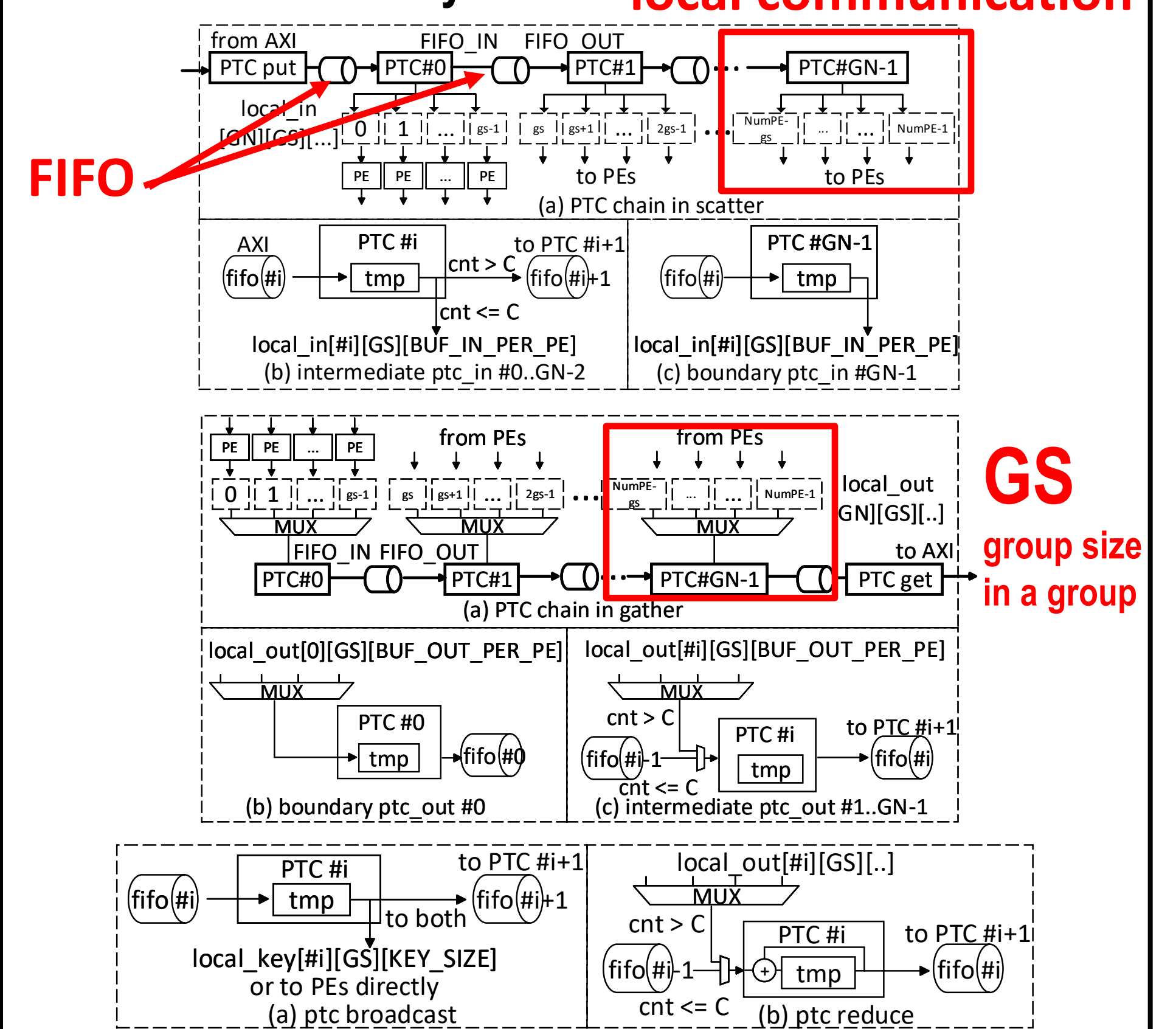


Broadcast and reduce are similar to the patterns of scatter and gather;

Pipelined Transfer Controller (PTC) in Latte

PTCs are chained in linear fashion through FIFOs.

Each PTC connects to a local set of buffers to constrain wire delay. **local communication**



Latte in HLS

260 lines of code (LOC) to manually implement Latte in HLS, 10x more than baseline code

PTC in Scatter Code snippet

```

1 #include <hls_stream>
2 int local_in[GN][GS][BUF_IN_PER_PE]; // redef.
3 void PTC_load(
4     int local_in[GN][GS][BUF_IN_PER_PE], int* global_in) {
5     #pragma HLS dataflow
6     hls::stream<int> fifo[GN]; // FIFOs, in Fig. 7a
7     ptc_put(global_in, fifo[0]);
8     for (int i = 0; i < GN-2; i++) {
9         ptc_in(fifo[i], fifo[i+1], local_in[i], GN-1-1);
10        ptc_in(fifo[GN-1], local_in[GN-1]);
11    }
12    void ptc_put(int* global_in, stream<int> &fifo) {
13        for (int i=0; i<NumPE; i++) {
14            for (int j = 0; j < BUF_IN_PER_PE; j++) {
15                #pragma HLS pipeline
16                fifo << global_in[i*BUF_IN_PER_PE+j];
17            }
18        }
19    }
20    void ptc_in( // #0..GN-2 ptc_in, in Fig. 7b
21        stream<int> &fifo_in, stream<int> &fifo_out,
22        int local_in[GN][GS][BUF_IN_PER_PE], int todo) {
23        int i, j, k; int tmp;
24        for (i=0; i < GS; i++) { // to local first
25            for (j = 0; j < BUF_IN_PER_PE; j++) {
26                tmp = fifo_in.read();
27                local_in[i][j] = tmp;
28            }
29        }
30        for (k=0; k < todo; k++) // to next ptc
31            for (i=0; i < GS; i++) {
32                for (j = 0; j < BUF_IN_PER_PE; j++) {
33                    tmp = fifo_in.read();
34                    fifo_out.write(tmp);
35                }
36            }
37    }
38 }
    
```

Can we do better? Yes!

Latte Automation

◆ Latte provides semiautomatic framework

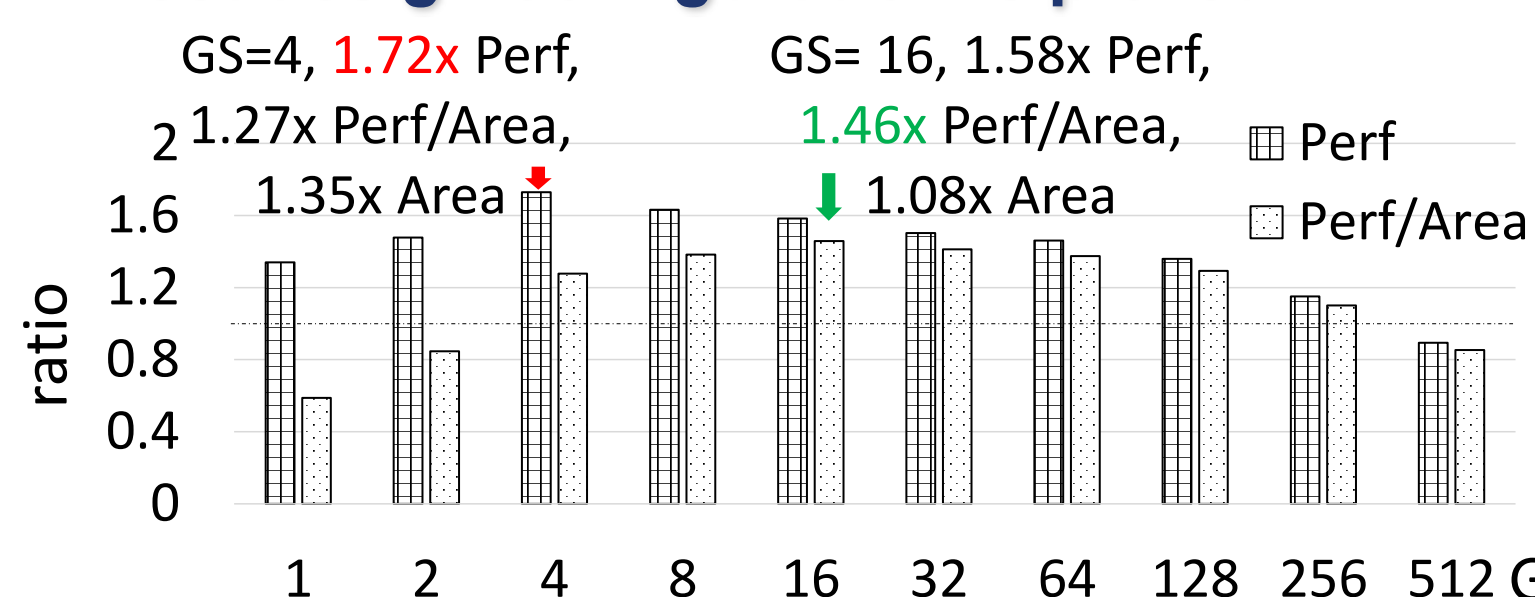
- Users insert simple Latte pragmas into user-written HLS kernel
- 260 -> 10 LOC, pragma indicates 1. buffers
- 2. pattern to be optimized

pattern variable that needs optimization

```

1 #pragma latte scatter var="local_in_B"
2 int local_in_B[NumPE][BUF_IN_PER_PE];
    
```

- Different PTC group sizes (#GS) are launched
- Best design configuration is picked



Latte is “Light Roast”

◆ Latte improves freq. with negligible overhead

- 1.50x with 3.2% LUT, 5.1% FF on average
- 2.66x with 2.7% LUT, 5.1% FF at max
- helps greatly in frequency degradation

linear layout of PTCs in gather pattern in FFT with 64 PEs and 16 PTCs

Baseline vs Latte

Bench.	type	N / GS	LUT	FF	DSP	BRAM	Freq.
AES	ori.	320 /	80.4%	17.3%	0.1%	76.3%	127
	latte	732 /	1.017	1.009	1	165.130x	57
FFT	ori.	64 /	80.5%	23.2%	88.9%	78.5%	132
	latte	744 /	1.027	1.056	1	207.158x	174
GEMM	ori.	512 /	87.8%	29.6%	71.1%	69.7%	141
	latte	716 /	1.044	0.962	1	192.120x	152
KMP	ori.	96 /	5.0%	3.0%	0.2%	52.3%	174
	latte	724 /	1.045	1.174	1	198.154x	174
NW	ori.	160 /	65.1%	50.7%	0.0%	78.2%	174
	latte	780 /	0.995	0.997	1	177.102x	174
SPMV	ori.	48 /	19.1%	11.9%	18.9%	93.2%	160
	latte	76 /	1.029	1.037	1	192.120x	174
STENCIL	ori.	64 /	12.9%	10.9%	48.1%	87.1%	141
	latte	716 /	1.044	1.139	1	198.154x	152
VITERBI	ori.	192 /	72.0%	25.7%	10.8%	39.3%	155
	latte	712 /	1.008	1.031	1	168.108x	152
Average	ori.	- /	NA	NA	NA	NA	120
	latte	- /	1.032	1.051	1	181.150x	152

