

FPGA-based Simulated Bifurcation Machine

Kosuke Tatsumura, Alexander R. Dixon, and Hayato Goto

FPL2019 @ Barcelona
2019.09.09

Contents

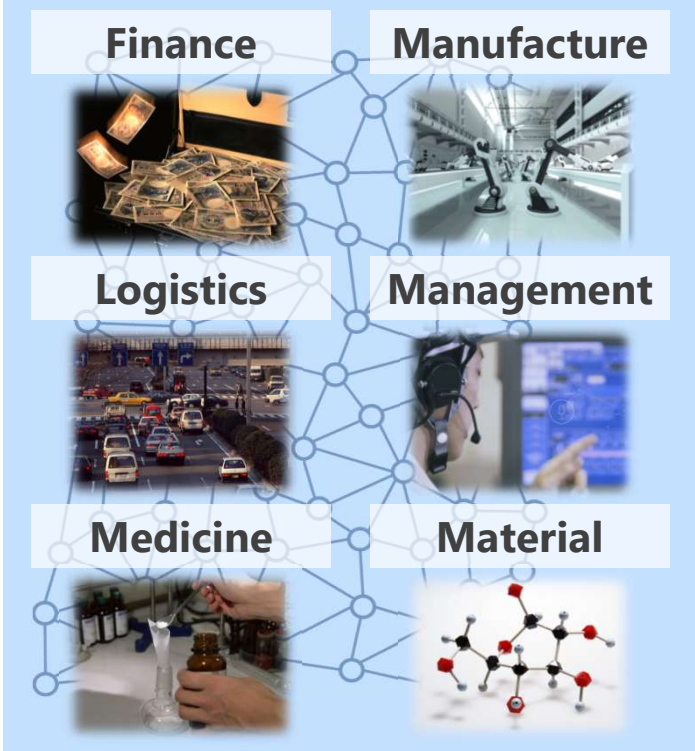
An ultra-fast solver for combinatorial optimization problems

- 01 Background & objective:
Accelerator for Simulated Bifurcation
- 02 Design
- 03 Implementation & Evaluation
- 04 Discussion: practicality

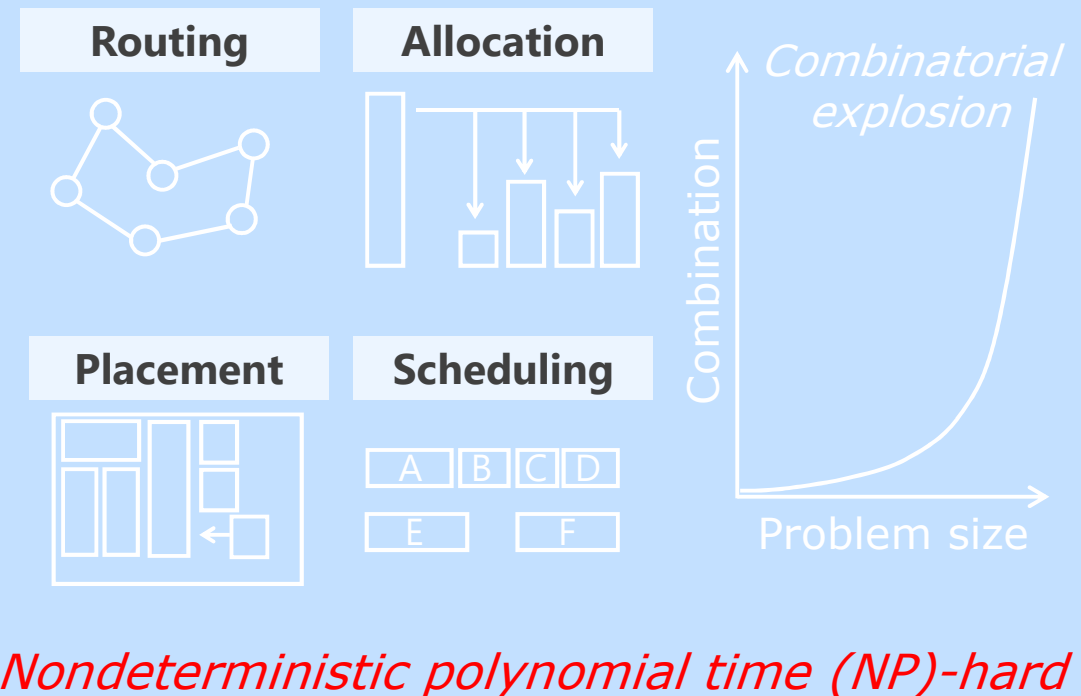
Combinatorial optimization problems

Economically valuable but computationally hard

Enhancing productivity



Combinatorial optimization problem



Standard approach: Simulated annealing (SA)

New approach: Ising machine

1. Map to an Ising problem,
2. Special-purpose machine solves it quickly

Problem

NP-hard

Combinatorial
optimization
problem

mapping

Ising problem

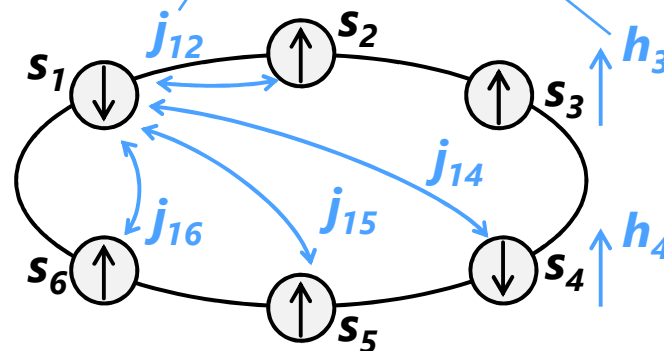
*NP-hard &
NP-complete*

$$J = \begin{bmatrix} 0 & j_{12} & j_{13} \\ j_{21} & 0 & j_{23} \\ j_{31} & j_{32} & 0 \end{bmatrix} \quad h = \begin{bmatrix} h_1 \\ h_2 \\ h_3 \end{bmatrix}$$

coupling

bias

input



spin: binary variable

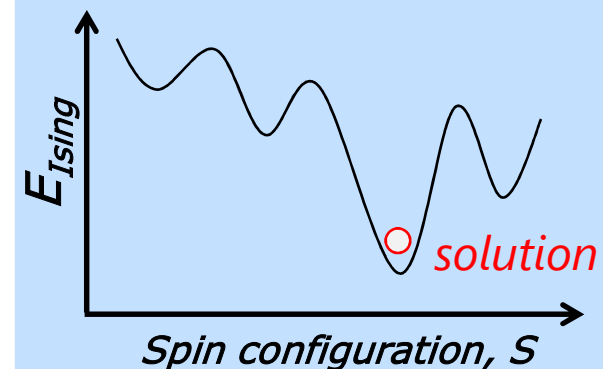
Ising machine

Special-purpose

search for ground-state \mathbf{s}
minimizing E

Ising Energy

$$E = - \sum j_{ij} s_i s_j + \sum h_i s_i$$



Ising machines

Very competitive, CIM deserves attention

D-Wave Sys.*¹
2011-

**Quantum
Annealer**



HITACHI *²
2015-

**CMOS annealing
machine**



FUJITSU *³
2016-

Digital annealer

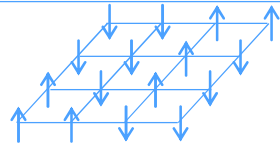


NTT/Stanford/U-Tokyo *⁴
2016-

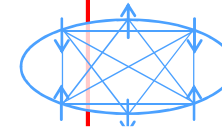
**Coherent Ising
machine (CIM)**



**Conne-
ctivity**



Locally connected
(sparse)



Fully connected
(dense)

Principle

Quantum
annealing

Simulated annealing (SA)

CIM
(optical bifurcation)

*1 <https://www.dwavesys.com/d-wave-two-system>

*2 <https://www.hitachi.co.jp/New/cnews/month/2019/02/0219.html>

*3 <https://www.fujitsu.com/global/about/resources/news/press-releases/2018/0515-01.html>

*4 <https://www.ntt.co.jp/news2017/1711e/171120a.html>

state-of-the-art

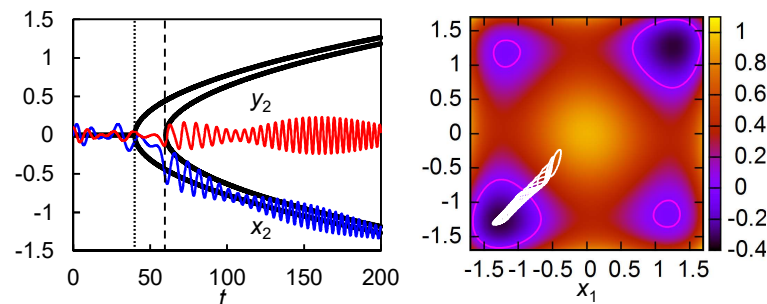
[R. Hamerly, *Sci. Adv.* eaau0823, '19]

A recently proposed, quantum-inspired algorithm

Simulated Bifurcation (SB)

[H. Goto, K. Tatsumura, A. R. Dixon, *Sci. Adv.* 5, eaau0823, '19]

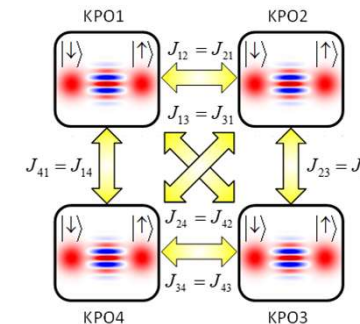
Simulated Bifurcation (SB)



Quantum Bifurcation (QB) machine

a quantum adiabatic optimization method

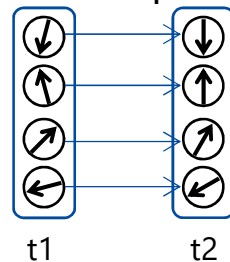
[H. Goto, *Sci. Rep.* 6, 21686, '16]



Derived as the
classical
counterpart

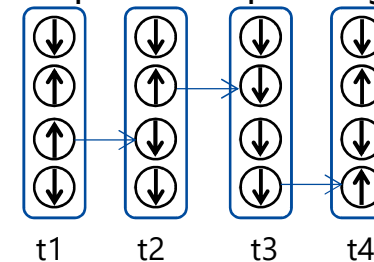
Plentiful parallelism

Parallel updating



Simulated Annealing (SA)

Sequential updating



→Substantial speedup by
massively parallel processing

Objective of this work

Designing & evaluating massively-parallel custom accelerator for SB

Motivation:

Be the world's fastest Ising machine

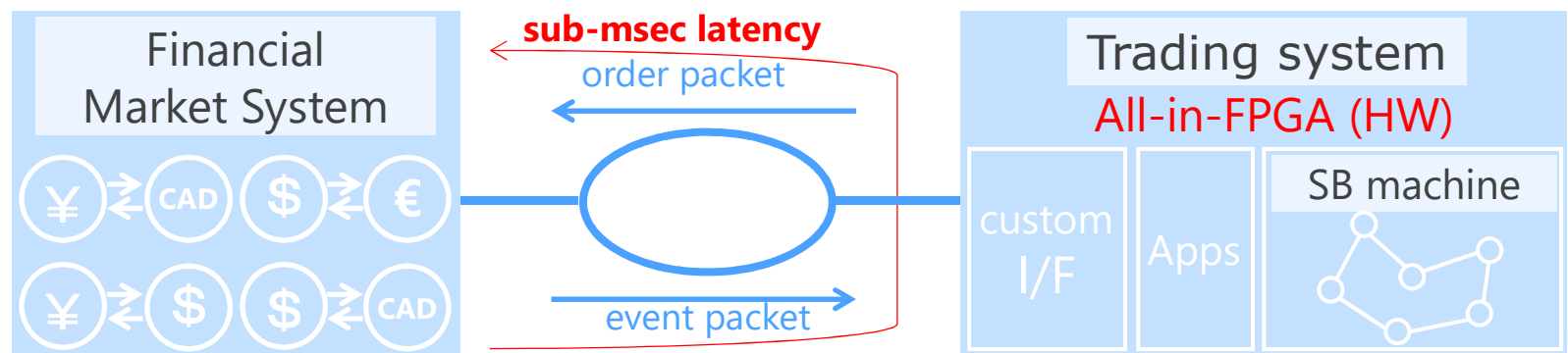
→ Pursue a custom circuit that can fully exploit the parallelism in SB

Meet the diversity of problems (size, bit precision, I/F)

→ Describe the design in an HLS to make it flexible & scalable

Toward “Intelligent real-time response systems”

→ Be a component of ALL-in-HW system enabling sub-msec latency



Contents

An ultra-fast solver for combinatorial optimization problems

- 01 Background & objective:
Accelerator for Simulated Bifurcation
- 02 Design
- 03 Implementation & Evaluation
- 04 Discussion: practicality

How it works: Simulated Bifurcation (SB)

***N*-body system dynamically searches for a good solution**

Movement of the system in *N*-dimensional space

Example: ***N*=2**

a single local minimum



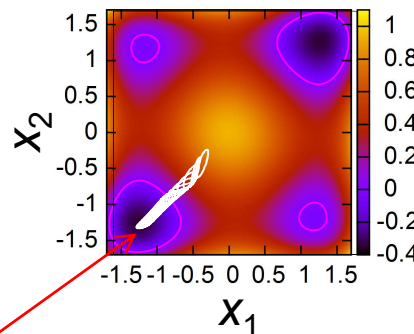
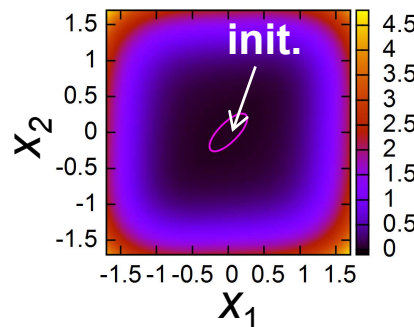
Bifurcation

(adiabatic process)



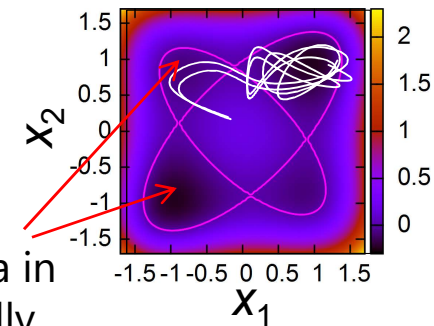
multiple local minima
(target cost function)

best solution
(-1,-1)

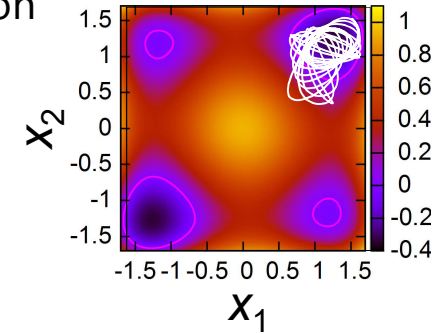


Adiabatic Search

chase one of the minima



Multiple minima in
the energetically
allowable region



Ergodic Search

find better one with higher probability

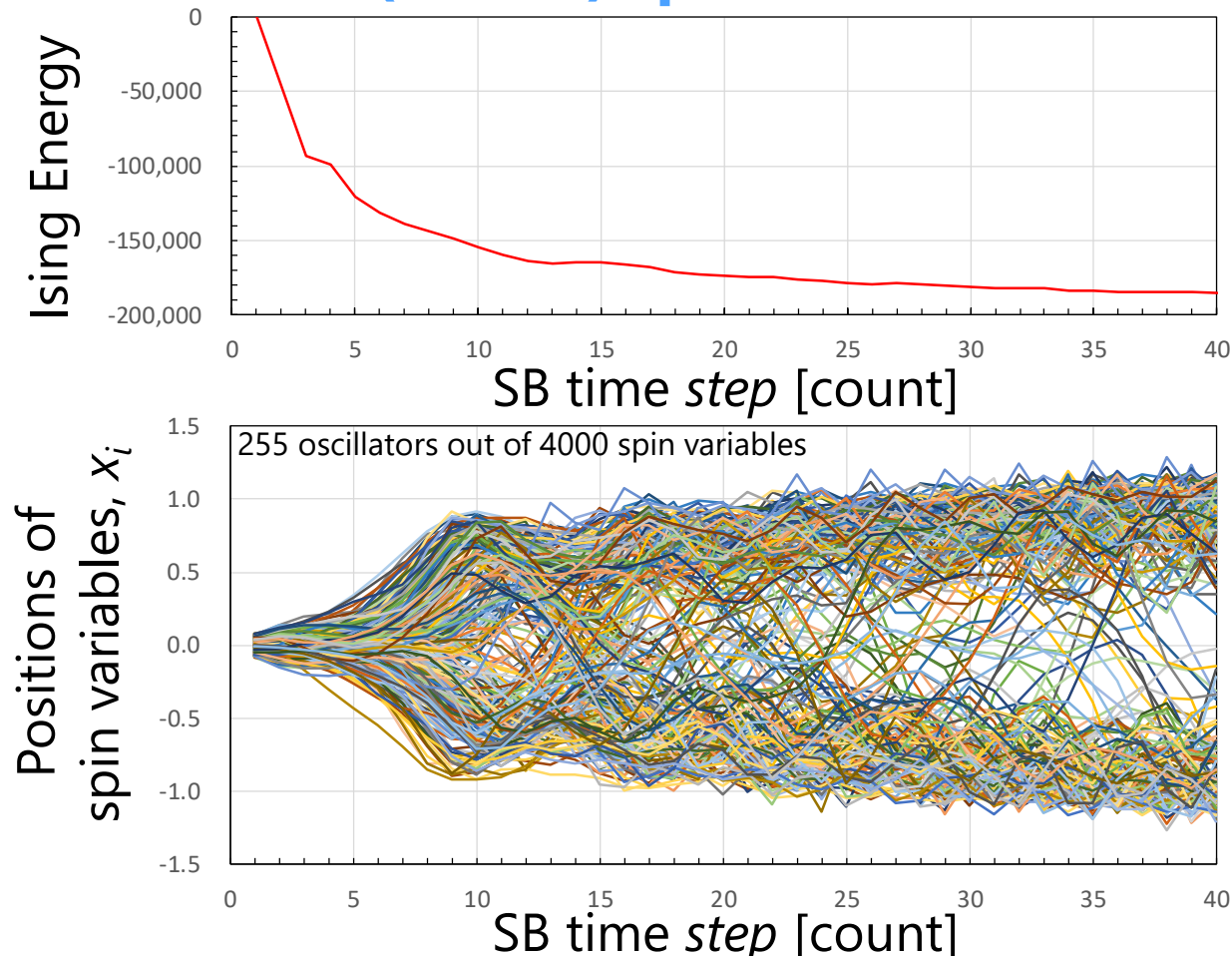
If *N* is large,

imagine the process of finding the brightest "star" among 2^N stars in the dark

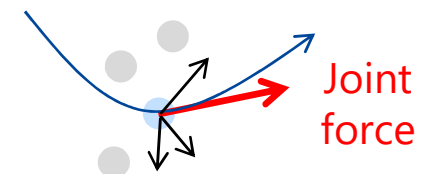
How it works: Simulated Bifurcation (SB)

Time evolution of N -body system

Movements of $N(=4000)$ spin-variables as a function of time



**better
solution**

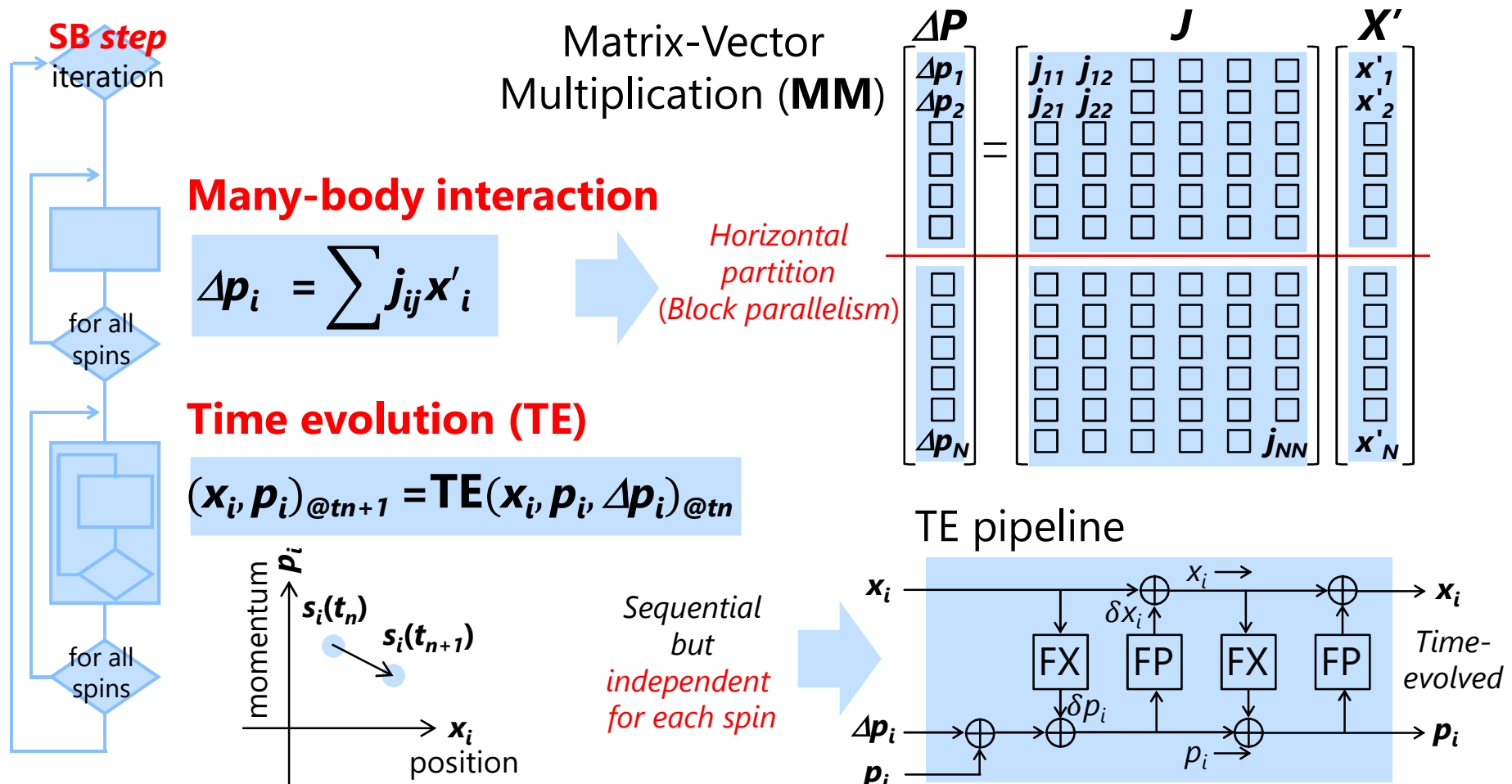


**Many-body
interactions**
depending on
all other spins

This is what FPGA-based SB machine computes

Algorithm of SB and it's parallelism

SB step: spin state at $t_{n+1} \leftarrow$ the previous state at t_n

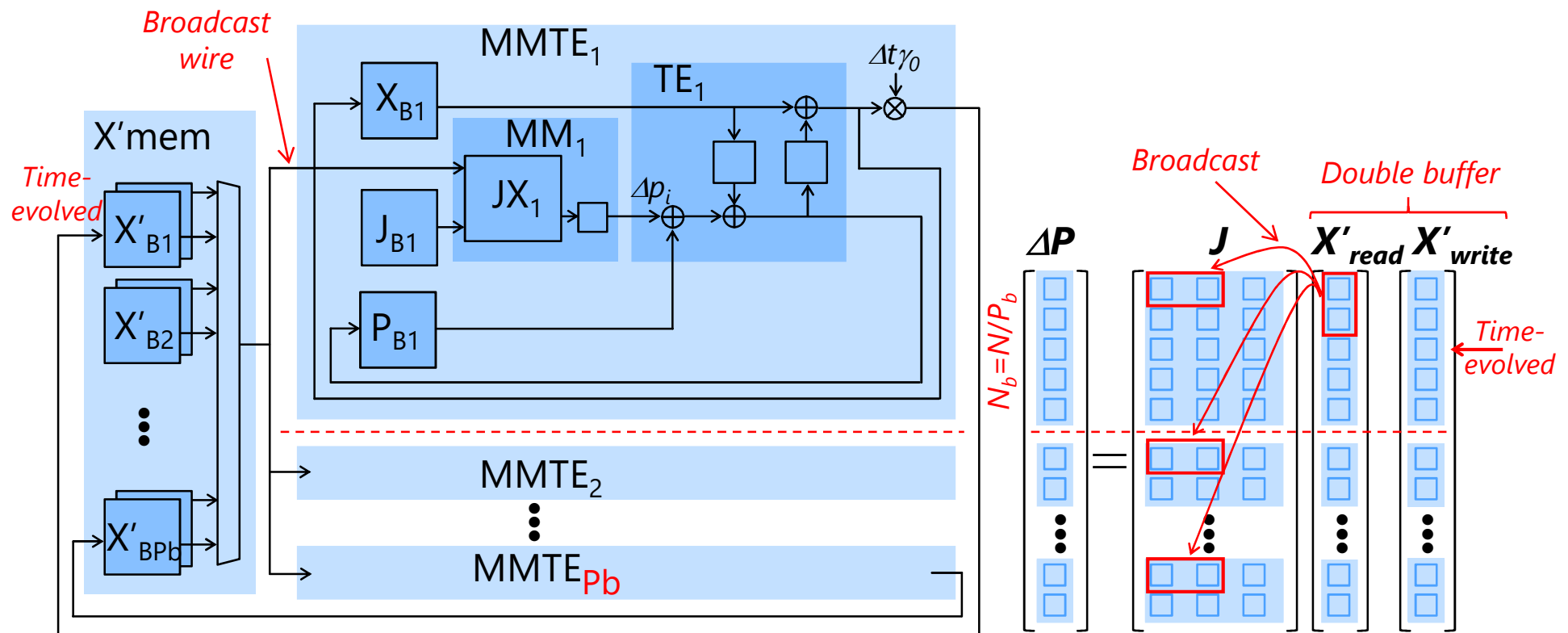


Top-level parallelism: **Simultaneous update of N spins is possible**

Design: Top-level circuit architecture

End-to-end HW implementation of SB

1. SB *step* iteration \rightarrow Circulative structure
2. Block parallelism (P_b) \rightarrow P_b number of MMTEs (responsible for $N_b = N/P_b$ spins)
3. Data dependency for \mathbf{x} \rightarrow *Double buffer* for X' mem



Design: the most computationally intensive part

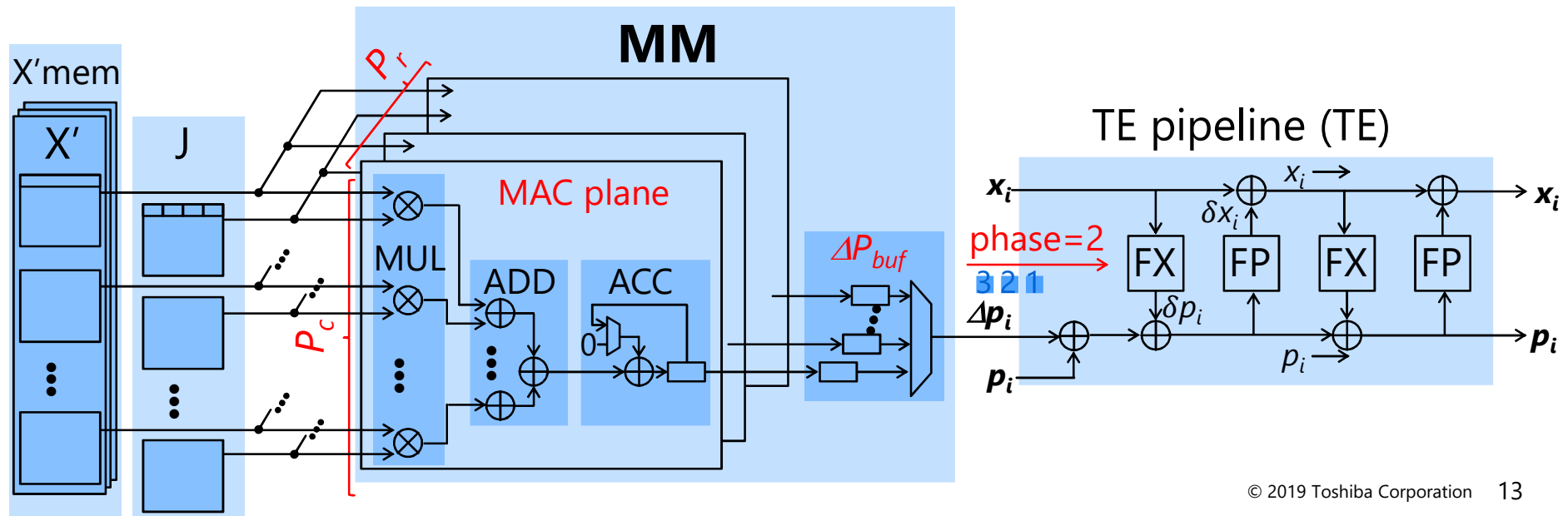
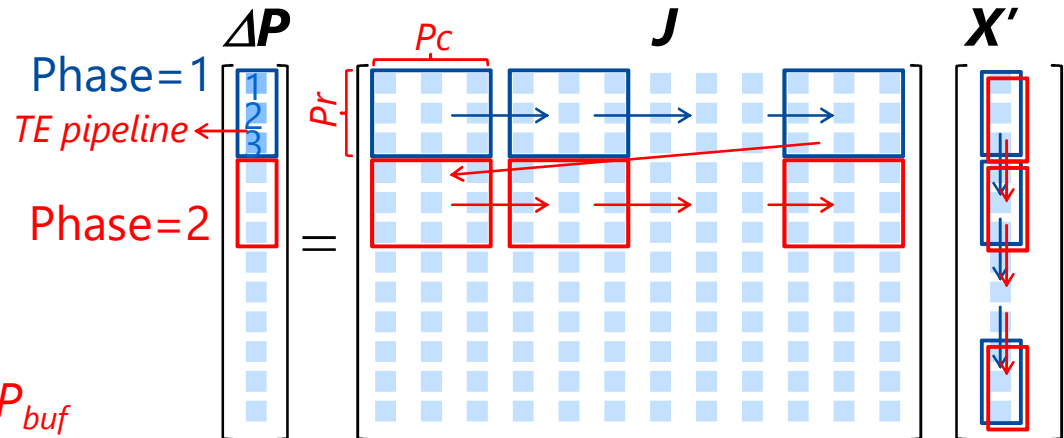
Matrix-vector multiplication module (MM)

1. Column/Row parallelism (P_c, P_r)

P_r number of P_c -input MAC planes

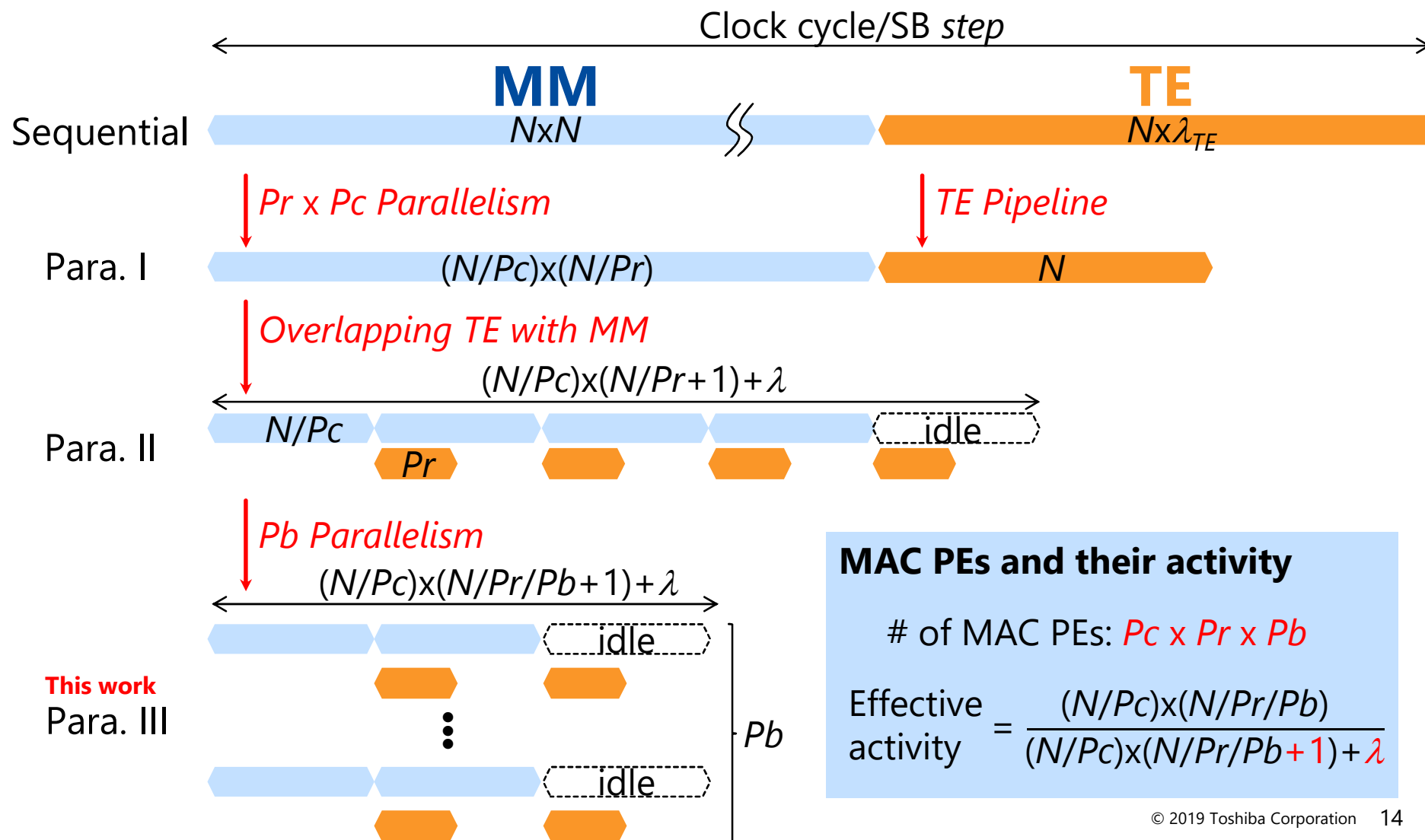
2. Overlapping TE with MM

Parallel-input/Sequential-output ΔP_{buf}



Timing Design: spatial and temporal parallelization

$P_c \times P_r \times P_b$ Speedup for MM & hiding TE part



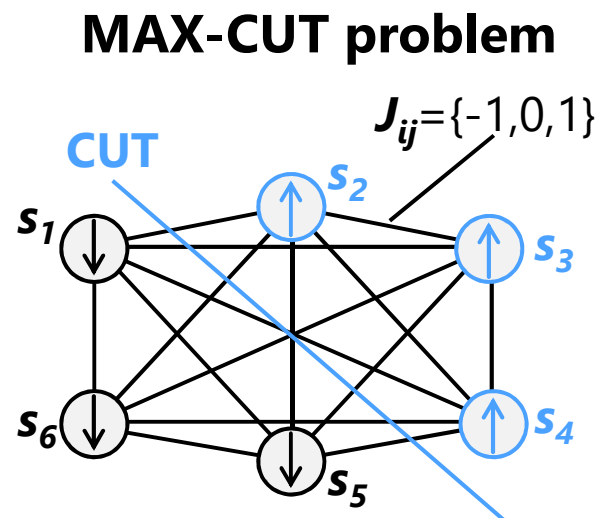
Contents

An ultra-fast solver for combinatorial optimization problems

- 01 Background & objective:
Accelerator for Simulated Bifurcation
- 02 Design
- 03 Implementation & Evaluation**
- 04 Discussion: practicality

Implementation for MAX-CUT benchmark problem

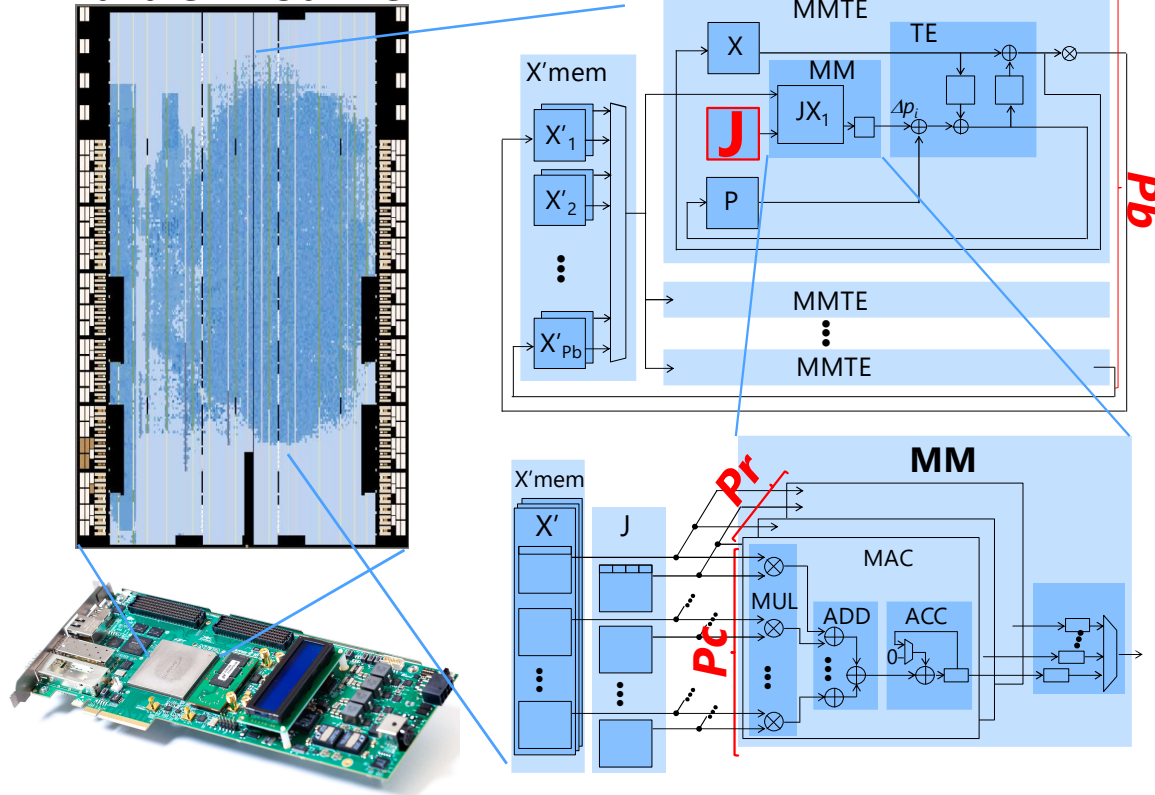
4,096-spin SB machine on Arria10 FPGA
Massively Parallel, High Utilization, 4X larger vs CIM



Implementation for MAX-CUT benchmark problem

4,096-spin SB machine on Arria10 FPGA Massively Parallel, High Utilization, 4X larger vs CIM

Arria10 GX1150 FPGA



Spin-Size: limited by **BRAM** (for J matrix)

MAX-CUT: $J_{ij} = \{-1, 0, 1\}$, $J_{4096 \times 4096} \rightarrow 16\text{Mbit}$ (1,024 BRAMs)

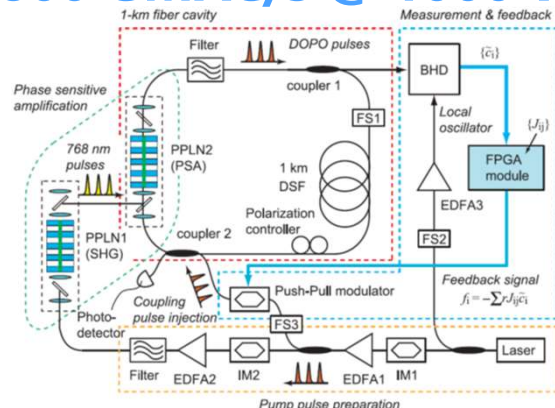
Parallelism: limited by **routing congestion**

	4,096-size
Architecture	
Pr/Pc/Pb	32/32/8
# of MAC PEs	8,192
Effective activity	92%
Resource	
ALM	40%
BRAM	56%
DSP	7%
System Clock	[MHz]
F_{sys}	269

Evaluation: FPGA-SB vs. CIM

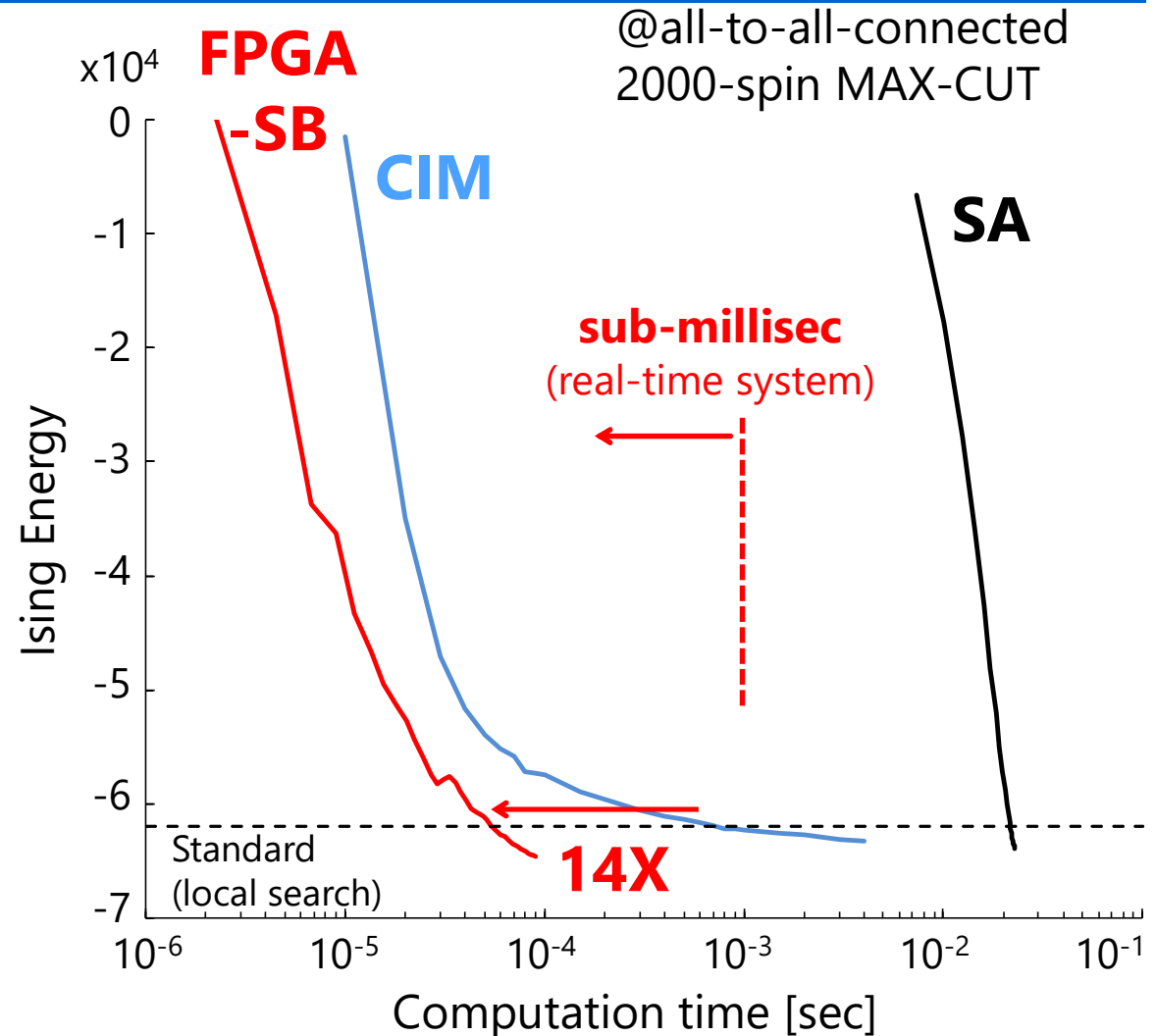
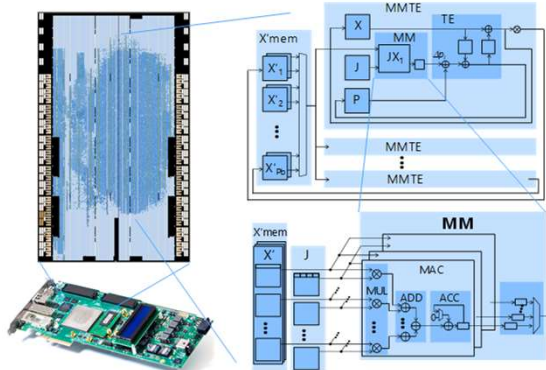
14X faster, 288X more energy efficient than CIM

Coherent Ising Machine
800 GMAC/s @ 1000 W



[T. Inagaki, Science 354, 603, '16]

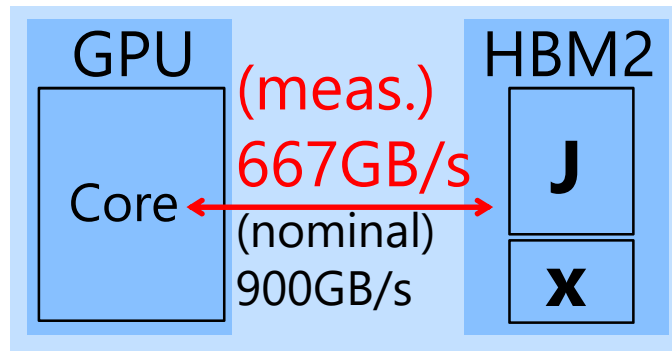
FPGA-SB
1,873 GMAC/s @ 49 W



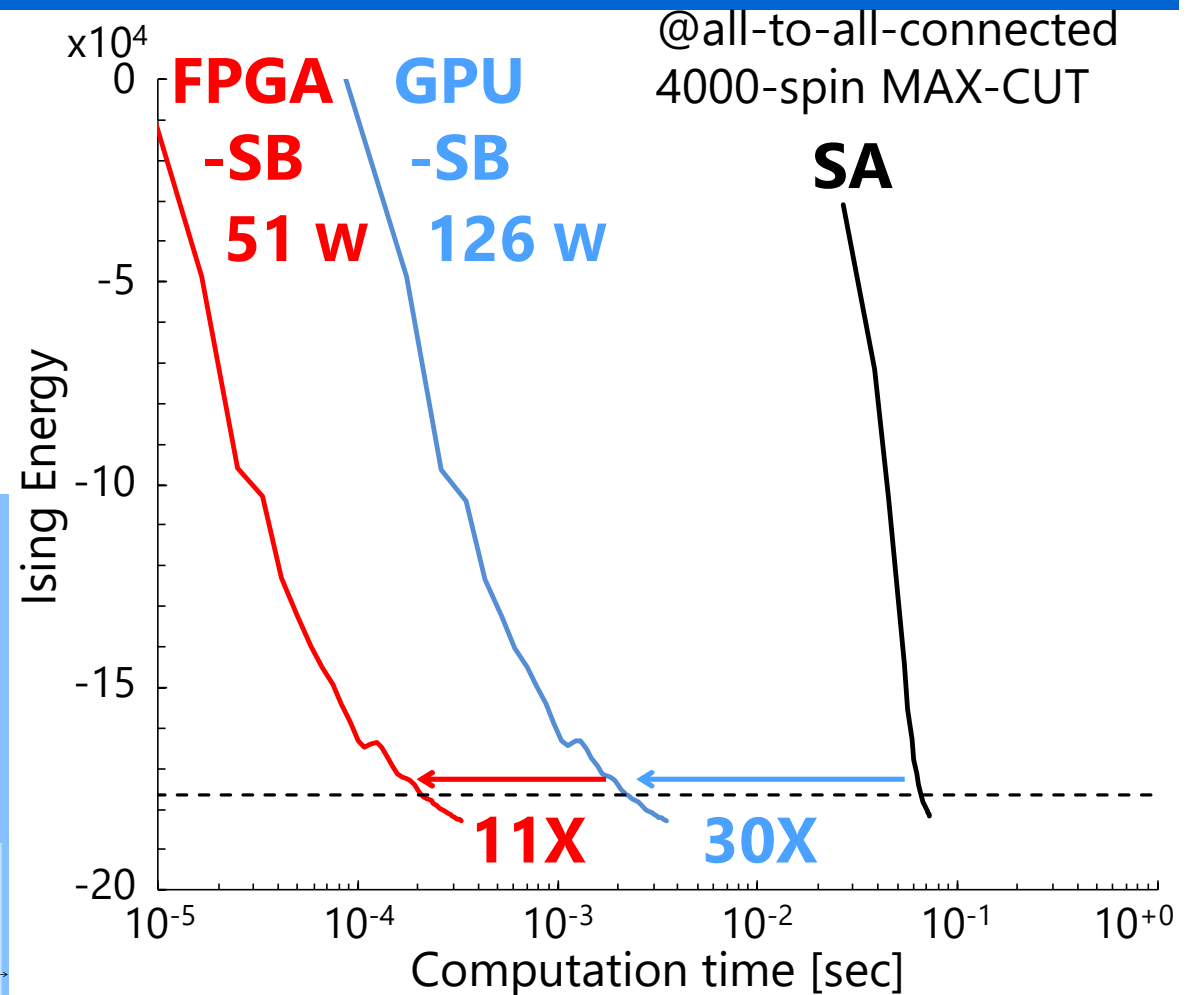
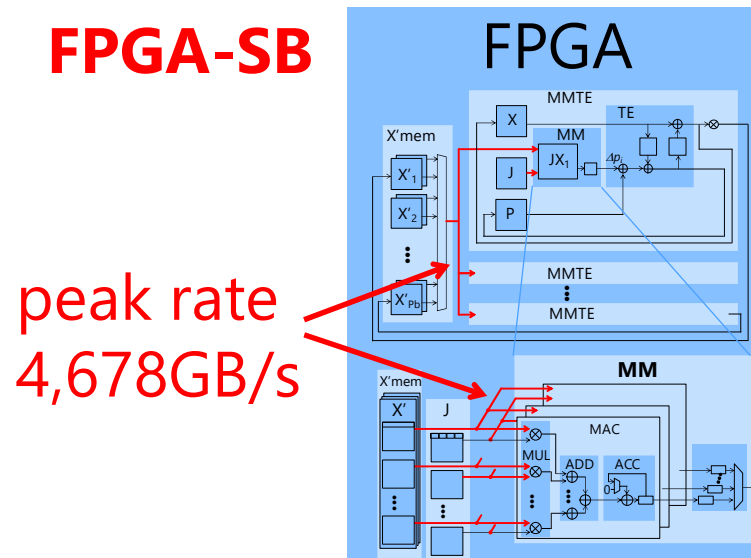
Evaluation: FPGA-SB vs. GPU-SB

FPGA is computation-bound, GPU memory-bound
-11X faster, 26X more energy efficient than GPU-SB

GPU-SB (Tesla V100)



FPGA-SB



Contents

An ultra-fast solver for combinatorial optimization problems

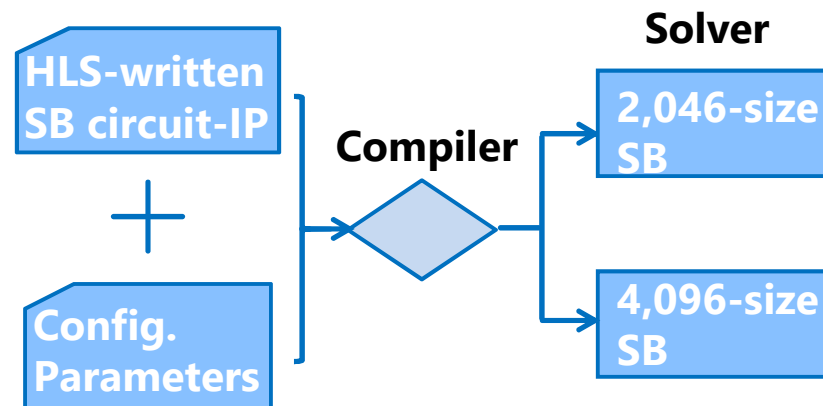
- 01 Background & objective:
Accelerator for Simulated Bifurcation
- 02 Design
- 03 Implementation & Evaluation
- 04 Discussion: practicality**

Hit record-high speed, meet diversity of problems

Write fully(bit-by-bit/clock-by-clock) customized HWs in an HLS language (OpenCL)

HLS: high-level-synthesis

Good size-scalability of our parameterized design



*all-to-all-connected MAX-CUT

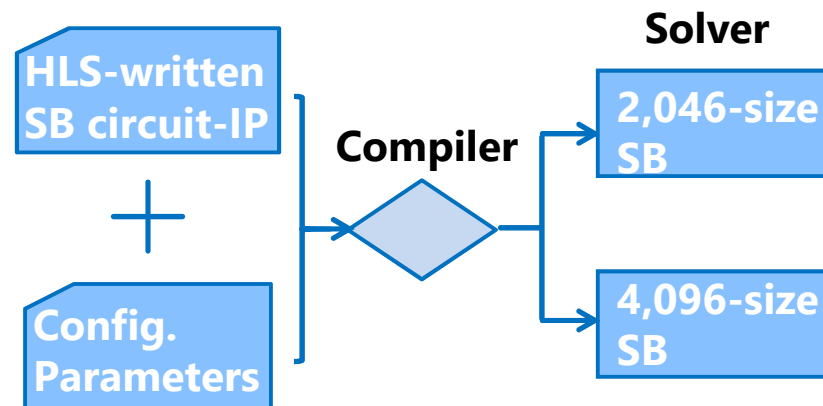
Problem	K2000*	K4000*	Ratio
Complexity per SB step	4M	16M	4X
Solver	2K-size SB	4K-size SB	Ratio
Time-to-solution [ms]	0.055	0.211	3.9X
Energy-for-solution [mJ]	2.7	10.8	4.0X

Hit record-high speed, meet diversity of problems

Write fully(bit-by-bit/clock-by-clock) customized HWs in an HLS language (OpenCL)

HLS: high-level-synthesis

Good size-scalability of our parameterized design



*all-to-all-connected MAX-CUT

Problem	K2000*	K4000*	Ratio
Complexity per SB step	4M	16M	4X
Solver	2K-size SB	4K-size SB	Ratio
Time-to-solution [ms]	0.055	0.211	3.9X
Energy-for-solution [mJ]	2.7	10.8	4.0X

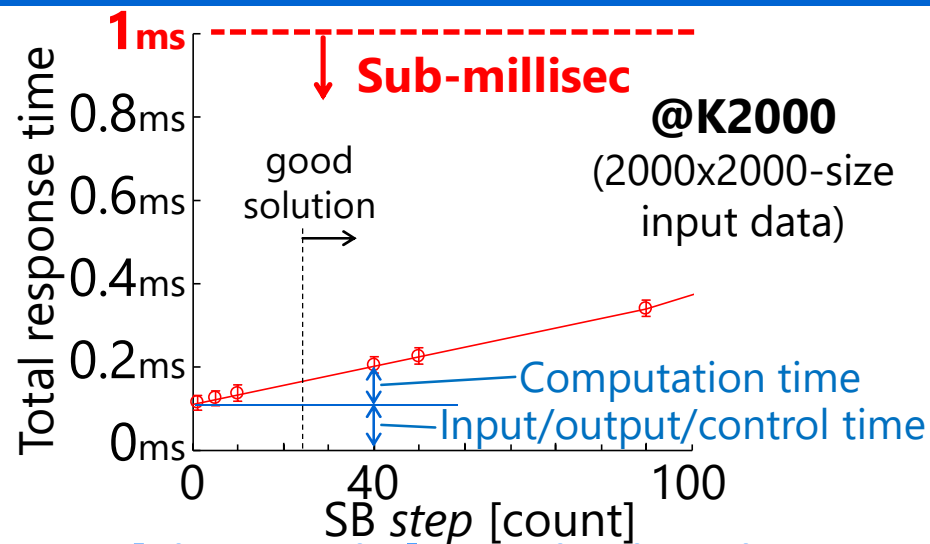
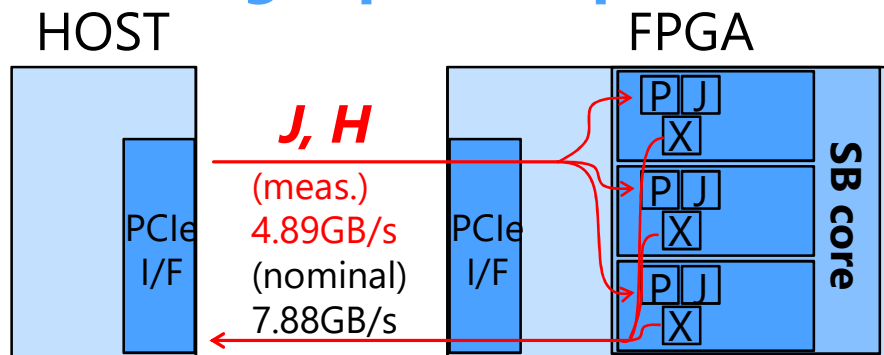
FPGA+HLS is a good choice for practical problems

Ising problems	J matrix	h vector	size
MAX-CUT	Ternary, {-1,0,1}	none	various
Arbitrage (path search)	Ternary, {-1,0,1}	FP32	various
Portfolio optimization	FP32	none	various
...	various	various	various

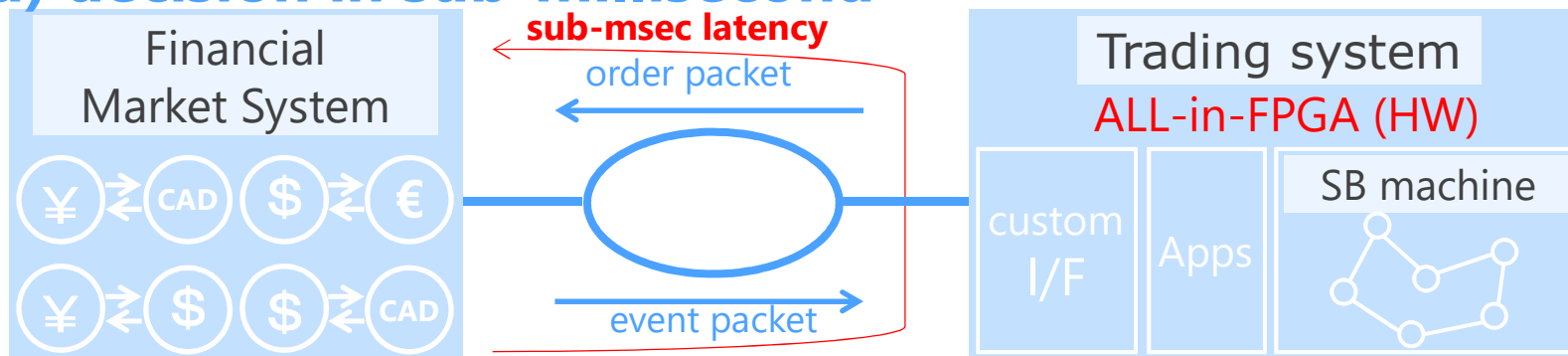
Toward “Intelligent real-time response systems”

Interface (I/F) is important as well;
Be flexible, wide, directly connected to external I/F

Total response time
including input/output/control



New concept: Making intelligent (combinatorial-optimization-based) decision in sub-millisecond



Summary

FPGA-based Simulated Bifurcation Machine

An ultra-fast solver for combinatorial optimization problems

Massively-parallel, fully-customized SB accelerator

- the world's fastest, most energy-efficient
(14X faster, 288X more energy efficient than CIM)

Very practical

- No refrigerator, no laser, but just an FPGA
- written in an HLS language, flexible & scalable to meet the diversity of problems

Be a key component for intelligent real-time response systems

- End-to-End HW implementation that can be directly connected to external systems, enabling *optimal* responses in sub-msec