

---

# **Integer Boltzmann Machines for Future Generation Digital Annealing Units**

Mohammad Bagherbeik<sup>1</sup>, Kouichi Kanda<sup>2</sup>,  
Hirotaka Tamura<sup>3</sup>, Ali Sheikholeslami<sup>1</sup>

<sup>1</sup>Dept. of Electrical & Computer Engineering, University of Toronto, Canada

<sup>2</sup>Fujitsu Limited, Kawasaki, Japan

<sup>3</sup>DXR Laboratories, Kawasaki, Japan

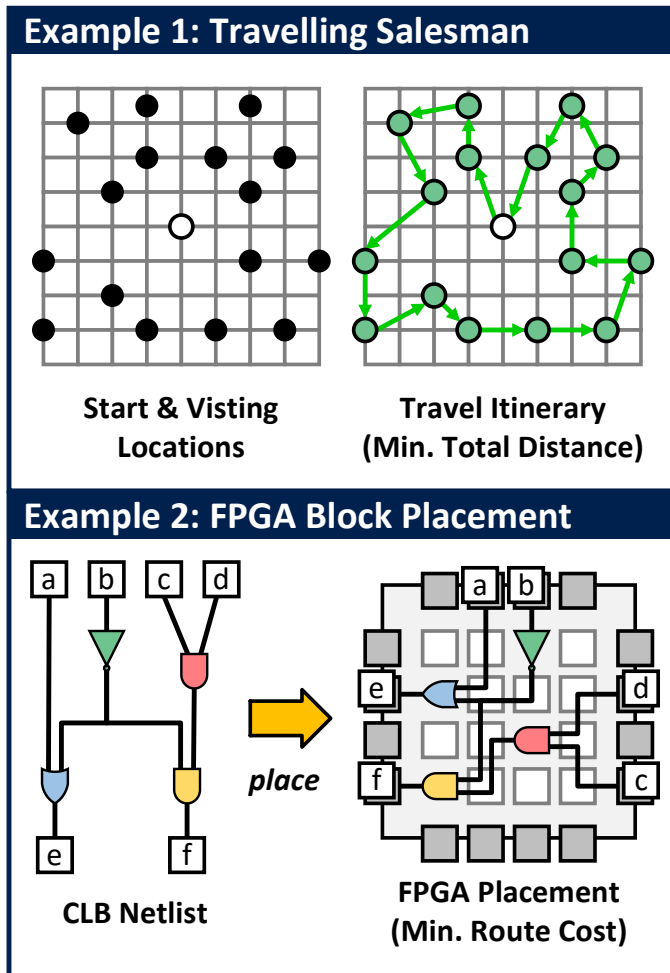
Adiabatic Quantum Computing Conference 2021

# Outline

---

- **Motivation and Background**
  - Quadratic Assignment Problems
  - Local Search Algorithms
  
- Integer Boltzmann Machine
  
- Parallel Tempering Solver Design
  - Algorithm and SIMD
  - Load Balancing
  
- Benchmark Results
  
- Conclusion

# Integer Assignment Problems

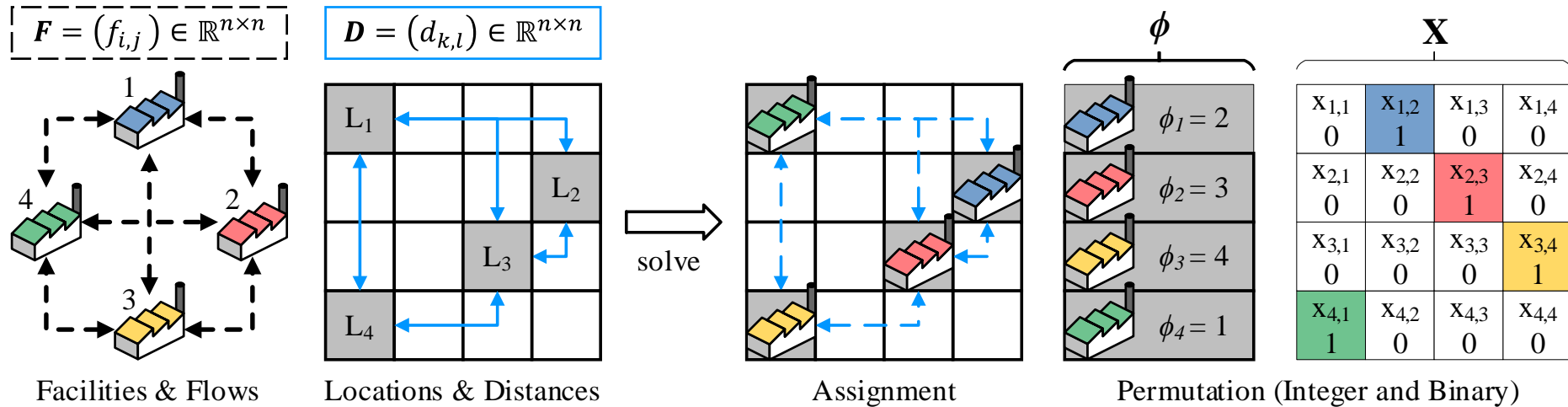


- Class of NP-hard combinatorial optimization problems with integer variables
  - Quadratic Assignment Problem (QAP) [1]
- Simple formulation but difficult to solve even with modern CPU and GPU hardware

## Goal

- Create an integer Boltzmann Machine variant to accelerate search for QAP
- Presentation covers formulation for symmetric, bias-less QAPs but concepts are directly transferrable to asymmetric + biased problems

# QAP Formulation



- Given  $n$  facilities and  $n$  locations
- Find a permutation  $\phi = (\phi_i) \in \Phi$ , where  $\Phi$  is set of all possible permutations
  - $\phi$  assigns each facility  $i$ , to a unique location  $\phi_i$  such that the total cost of the assignment (1) is minimized:

$$C(\phi) = \sum_{i=1}^n \sum_{j=1}^n f_{i,j} d_{\phi_i, \phi_j} \quad (1)$$

$$x_{i,j} = \begin{cases} 1, & \text{if } \phi_i = j \\ 0, & \text{otherwise} \end{cases}, \mathbf{X} = (x_{i,j}) \in \{0,1\}^{n \times n} \quad (2)$$

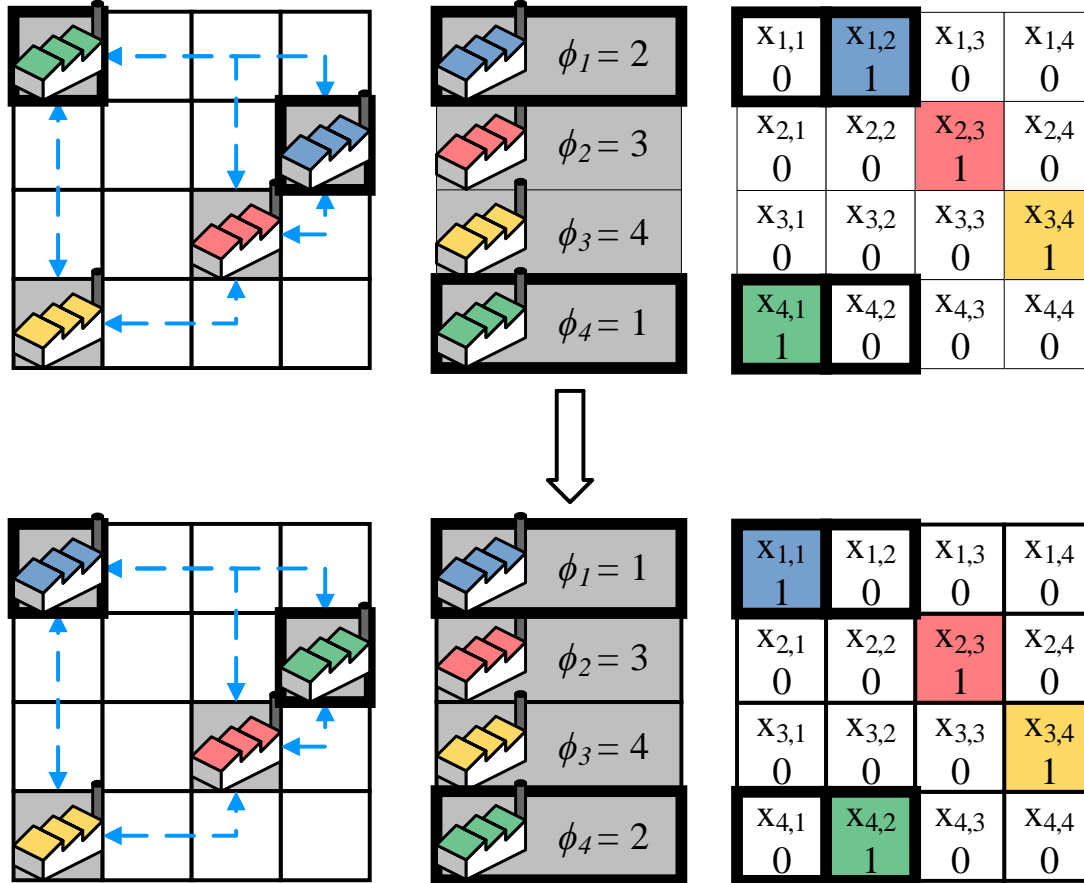
# Local Search

---

1. Propose new states as perturbations of the current state
  - Exchange locations of 2 facilities
2. Calculate difference in cost function ( $\Delta C$ ) between previous and proposed state
3. Decide whether to accept proposed state based on some criterion
  - Stochastic Approach (e.g. Simulated Annealing):  
Use a varying temperature variable ( $T$ ) that is used to inject randomness into acceptance of moves via Metropolis Hastings acceptance probability generated using (3)
    - Proposal Acceptance Rate ( $PAR$ )
    - $T \rightarrow \infty, PAR \rightarrow 100\%$ , all proposals are accepted regardless of  $\Delta C$  value
    - $T \rightarrow 0, PAR \rightarrow 0\%$ , search turns greedy

$$P_{acc} = \min\left\{1, \exp\left(-\frac{\Delta C}{T}\right)\right\} \quad (3)$$

# Local Search: QAP Exchange



- Exchange locations of two facilities  $a, b$
- Change in cost calculated via (4)
- Implemented in software via simple loop(s) in  $\mathcal{O}(n)$  time

$$\Delta C_{ex} = \sum_{i=1, i \neq a, b}^n 2(f_{b,i} - f_{a,i})(d_{\phi_a, \phi_i} - d_{\phi_b, \phi_i}) \quad (4)$$

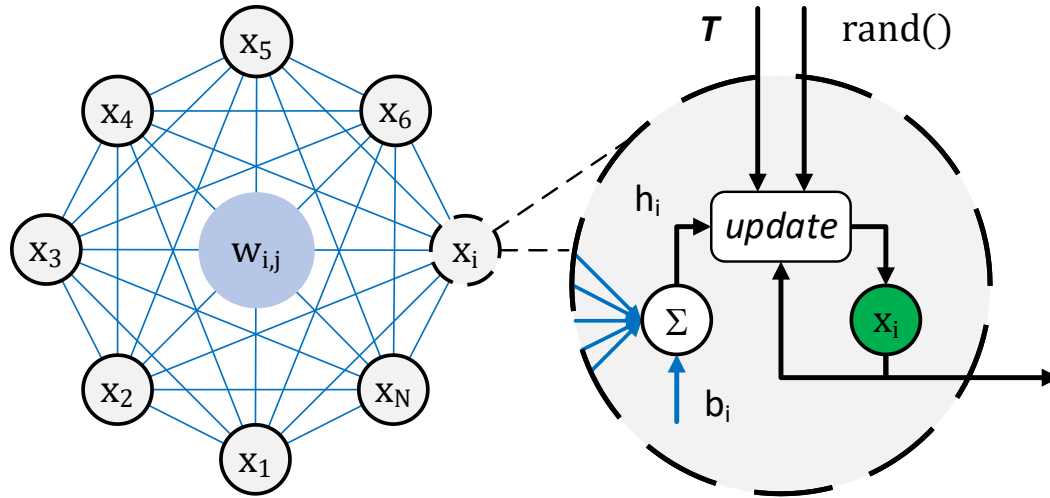
# Outline

---

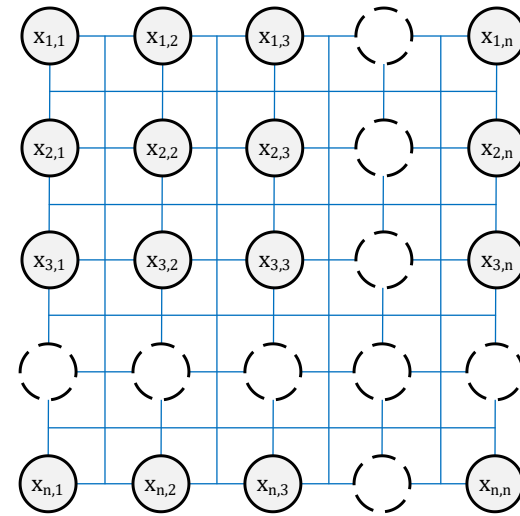
- Motivation and Background
  - Quadratic Assignment Problems
  - Local Search Algorithms
  
- **Integer Boltzmann Machine**
  
- Parallel Tempering Solver Design
  - Algorithm and SIMD
  - Load Balancing
  
- Benchmark Results
  
- Conclusion

# Boltzmann Machines

1D -  $N$  neurons



2D -  $n \times n$  neurons



State Cost/ Energy	$C(\mathbf{X}) = -\frac{1}{2} \sum_{i=1}^N \sum_{j=1}^N w_{i,j} x_i x_j - \sum_{i=1}^N b_i x_i$	(5)	$C(\mathbf{X}) = -\frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \sum_{k=1}^n \sum_{l=1}^n w_{i,j;k,l} x_{i,j} x_{k,l} - \sum_{i=1}^n \sum_{j=1}^n b_{i,j} x_{i,j}$	(7)
Neuron Local Fields	$h_i(\mathbf{X}) = \sum_{j=1}^N w_{i,j} x_j + b_i$	(6)	$h_{i,j}(\mathbf{X}) = \sum_{k=1}^n \sum_{l=1}^n w_{i,j;k,l} x_{k,l} + b_{i,j}$	(8)



# BM vs Integer BM for QAP [2,3]

	Data Storage	Searching the State-Space: Exchange 2 Facilities (a↔b)
<b>BM</b>	<p>n locations</p> <p>n facilities</p> <p><b>State</b> n<sup>2</sup> binary bits 64kbits at n=256</p> <p><b>W matrix:</b> Fully Connected Bits 1.6G floats at n=256</p>	<p>S1      S2      S3      S4      S5</p> <p><b>Infeasible States</b></p> <p>4 Sequential Spin-Flips</p> <p>Cost(X)</p> <p>Search State-Space</p>
<b>Int BM</b>	<p><math>\phi</math></p> <p><b>State</b> n integers 2,048bits at n=256</p> <p><b>F, D Matrices</b> 2 x 64k floats at n=256</p>	$\Delta C_{ex}(a \leftrightarrow b) = 2(h_{a,\phi_b} + h_{b,\phi_a} - h_{a,\phi_a} - h_{b,\phi_b}) + 4f_{a,b}d_{\phi_a\phi_b}$ <p>S1      S5</p> <p>1 Exchange Operation</p> <p>Cost(<math>\phi</math>)</p> <p>Search State-Space</p>

- ❑ 2D BM requires penalty terms to enforce one-hot constraints on rows/columns of  $X$
- ❑ Using an Integer BM with intrinsic exchange moves (4-bit-flip equivalent) simplifies problem formulation and speeds up search, significantly reduces memory reqs.

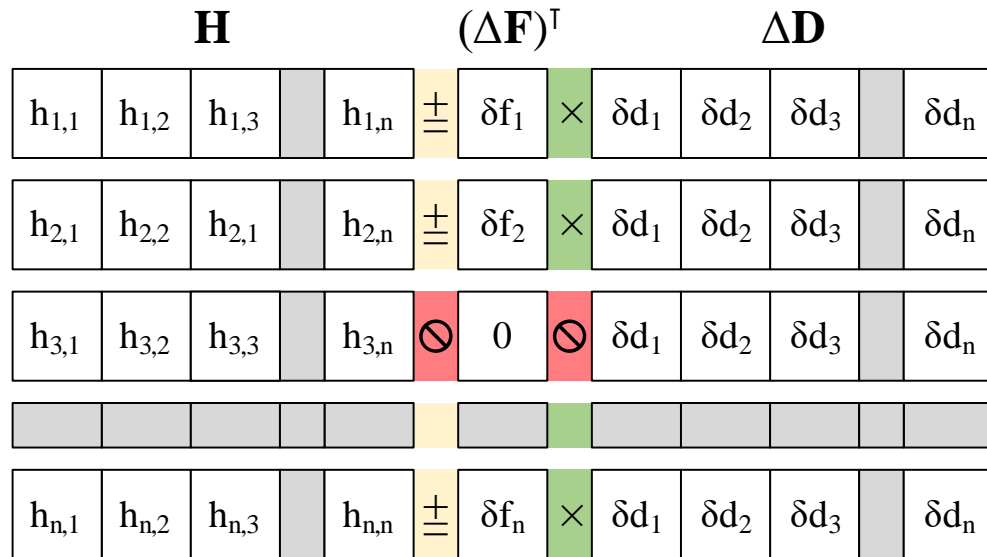
# Int BM QAP Formulation

Notes	Formulation	#
<ul style="list-style-type: none"> <li>Remove conditions on the summation loop iterator</li> </ul>	$\Delta C_{ex} = \sum_{i=1, i \neq a, b}^n 2(f_{b,i} - f_{a,i})(d_{\phi_a, \phi_i} - d_{\phi_b, \phi_i})$	(4)
<ul style="list-style-type: none"> <li>Introduces extra terms within the sum</li> <li>Compensate via additional terms</li> </ul>	$\Delta C_{ex} = \sum_{i=1}^n 2(f_{b,i} - f_{a,i})(d_{\phi_a, \phi_i} - d_{\phi_b, \phi_i}) + 4f_{a,b}d_{\psi_a\psi_b}$	(9)
<ul style="list-style-type: none"> <li>Create vector of flow row differences</li> </ul>	$\Delta_a^b \mathbf{F} = \mathbf{F}_{b,*} - \mathbf{F}_{a,*}$	(10)
<ul style="list-style-type: none"> <li>Create vector of permutation ordered distance row differences</li> </ul>	$\Delta_{\phi_b}^{\phi_a} \mathbf{D}^{\mathbf{X}} = (\mathbf{D}_{\phi_a,*} - \mathbf{D}_{\phi_b,*}) \mathbf{X}^{\top}$	(11)
<ul style="list-style-type: none"> <li>Reformulate using a dot product</li> </ul>	$\Delta C_{ex} = 2\Delta_a^b \mathbf{F} \cdot \Delta_{\phi_b}^{\phi_a} \mathbf{D}^{\mathbf{X}} + 4f_{a,b}d_{\phi_a\phi_b}$	(12)

# Integer BM: Local Fields

Notes	Formulation	#																																														
<ul style="list-style-type: none"> <li>Generate cache values at start of search in <math>\mathcal{O}(n^3)</math> time</li> </ul>	$\mathbf{H} = \mathbf{F}(\mathbf{D}^\top \mathbf{X}) = (h_{r,c}) \in \mathbb{R}^{n \times n}$	(13)																																														
<ul style="list-style-type: none"> <li>Elements in <math>\mathbf{H}</math> correspond to bits in <math>\mathbf{X}</math></li> <li>Element <math>h_{r,c}</math> is cost of assigning a facility <math>r</math> to location <math>c</math></li> </ul>	<div style="display: flex; justify-content: space-around; align-items: center;"> <table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr><th><math>\mathbf{X}</math></th><th>L<sub>1</sub></th><th>L<sub>2</sub></th><th>L<sub>3</sub></th><th>L<sub>4</sub></th></tr> </thead> <tbody> <tr><td>0</td><td>x<sub>1,1</sub></td><td>x<sub>1,2</sub></td><td>x<sub>1,3</sub></td><td>x<sub>1,4</sub></td></tr> <tr><td>0</td><td>x<sub>2,1</sub></td><td>x<sub>2,2</sub></td><td>x<sub>2,3</sub></td><td>x<sub>2,4</sub></td></tr> <tr><td>0</td><td>x<sub>3,1</sub></td><td>x<sub>3,2</sub></td><td>x<sub>3,3</sub></td><td>x<sub>3,4</sub></td></tr> <tr><td>1</td><td>x<sub>4,1</sub></td><td>x<sub>4,2</sub></td><td>x<sub>4,3</sub></td><td>x<sub>4,4</sub></td></tr> </tbody> </table> <table border="1" style="border-collapse: collapse; text-align: center;"> <thead> <tr><th><math>\mathbf{H}</math></th><th>L<sub>1</sub></th><th>L<sub>2</sub></th><th>L<sub>3</sub></th><th>L<sub>4</sub></th></tr> </thead> <tbody> <tr><td>h<sub>1,1</sub></td><td>h<sub>1,2</sub></td><td>h<sub>1,3</sub></td><td>h<sub>1,4</sub></td></tr> <tr><td>h<sub>2,1</sub></td><td>h<sub>2,2</sub></td><td>h<sub>2,3</sub></td><td>h<sub>2,4</sub></td></tr> <tr><td>h<sub>3,1</sub></td><td>h<sub>3,2</sub></td><td>h<sub>3,3</sub></td><td>h<sub>3,4</sub></td></tr> <tr><td>h<sub>4,1</sub></td><td>h<sub>4,2</sub></td><td>h<sub>4,3</sub></td><td>h<sub>4,4</sub></td></tr> </tbody> </table> </div>	$\mathbf{X}$	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	0	x <sub>1,1</sub>	x <sub>1,2</sub>	x <sub>1,3</sub>	x <sub>1,4</sub>	0	x <sub>2,1</sub>	x <sub>2,2</sub>	x <sub>2,3</sub>	x <sub>2,4</sub>	0	x <sub>3,1</sub>	x <sub>3,2</sub>	x <sub>3,3</sub>	x <sub>3,4</sub>	1	x <sub>4,1</sub>	x <sub>4,2</sub>	x <sub>4,3</sub>	x <sub>4,4</sub>	$\mathbf{H}$	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>	h <sub>1,1</sub>	h <sub>1,2</sub>	h <sub>1,3</sub>	h <sub>1,4</sub>	h <sub>2,1</sub>	h <sub>2,2</sub>	h <sub>2,3</sub>	h <sub>2,4</sub>	h <sub>3,1</sub>	h <sub>3,2</sub>	h <sub>3,3</sub>	h <sub>3,4</sub>	h <sub>4,1</sub>	h <sub>4,2</sub>	h <sub>4,3</sub>	h <sub>4,4</sub>	
$\mathbf{X}$	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>																																												
0	x <sub>1,1</sub>	x <sub>1,2</sub>	x <sub>1,3</sub>	x <sub>1,4</sub>																																												
0	x <sub>2,1</sub>	x <sub>2,2</sub>	x <sub>2,3</sub>	x <sub>2,4</sub>																																												
0	x <sub>3,1</sub>	x <sub>3,2</sub>	x <sub>3,3</sub>	x <sub>3,4</sub>																																												
1	x <sub>4,1</sub>	x <sub>4,2</sub>	x <sub>4,3</sub>	x <sub>4,4</sub>																																												
$\mathbf{H}$	L <sub>1</sub>	L <sub>2</sub>	L <sub>3</sub>	L <sub>4</sub>																																												
h <sub>1,1</sub>	h <sub>1,2</sub>	h <sub>1,3</sub>	h <sub>1,4</sub>																																													
h <sub>2,1</sub>	h <sub>2,2</sub>	h <sub>2,3</sub>	h <sub>2,4</sub>																																													
h <sub>3,1</sub>	h <sub>3,2</sub>	h <sub>3,3</sub>	h <sub>3,4</sub>																																													
h <sub>4,1</sub>	h <sub>4,2</sub>	h <sub>4,3</sub>	h <sub>4,4</sub>																																													
<ul style="list-style-type: none"> <li>Reformulate dot product as sum of 4 <math>\mathbf{H}</math> elements taking <math>\mathcal{O}(1)</math> time</li> </ul>	$\Delta_a^b \mathbf{F} \cdot \Delta_{\phi_b}^{\phi_a} \mathbf{D}^{\mathbf{X}} = h_{a,\phi_b} + h_{b,\phi_a} - h_{a,\phi_a} - h_{b,\phi_b}$	(14)																																														
	$\Delta C_{ex} = 2\Delta_a^b \mathbf{F} \cdot \Delta_{\phi_b}^{\phi_a} \mathbf{D}^{\mathbf{X}} + 4f_{a,b} d_{\phi_a \phi_b} \quad \rightarrow \quad \Delta C_{ex} = 2(h_{a,\phi_b} + h_{b,\phi_a} - h_{a,\phi_a} - h_{b,\phi_b}) + 4f_{a,b} d_{\phi_a \phi_b}$	(15)																																														

# Int BM: Local Field Updates



- Update local fields 1 row at a time using floating point fused-multiply-adds
  - $\mathcal{O}(n^2)$
- Can skip rows corresponding to 0 elements in  $\Delta\mathbf{F}$
- Significant speed-ups if  $\mathbf{F}$  is sparse or has a particular pattern that leads to sparse  $\Delta\mathbf{F}$  vectors

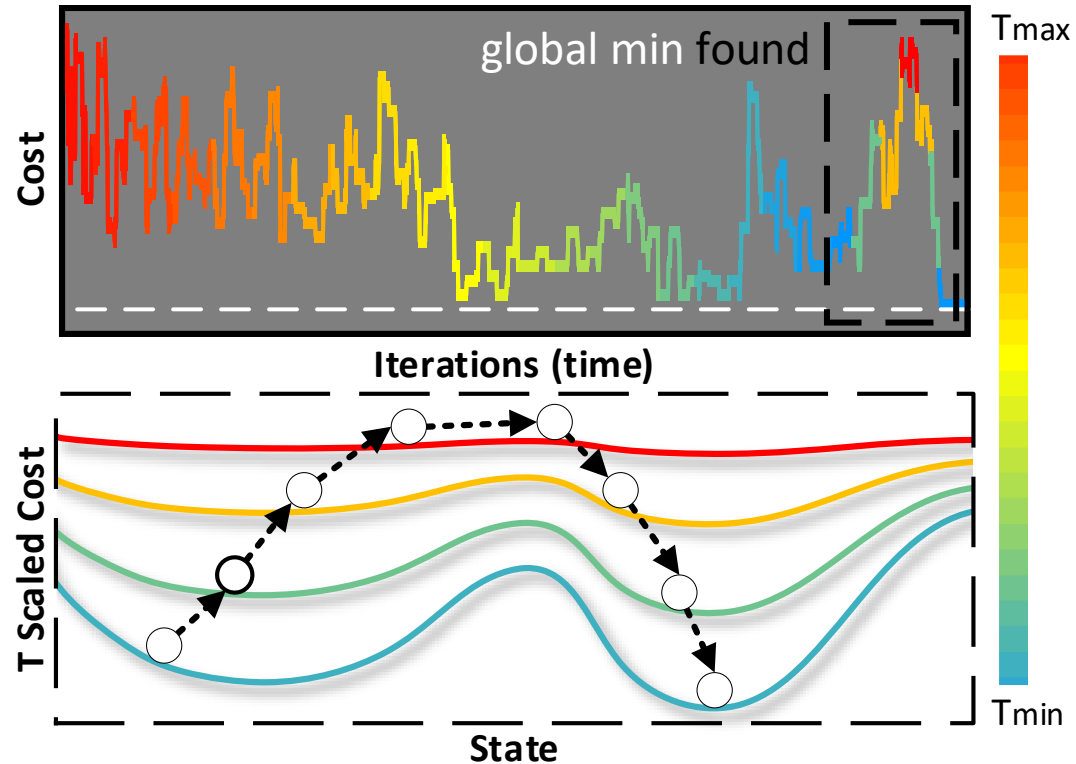
$$\mathbf{H} += (\Delta_a^b \mathbf{F})^\top \otimes \Delta_{\phi_b}^{\phi_a} \mathbf{D} \quad (16)$$

# Outline

---

- Motivation and Background
  - Quadratic Assignment Problems
  - Local Search Algorithms
  
- Integer Boltzmann Machine
  
- **Parallel Tempering Solver Design**
  - Algorithm and SIMD
  - Load Balancing
  
- Benchmark Results
  
- Conclusion

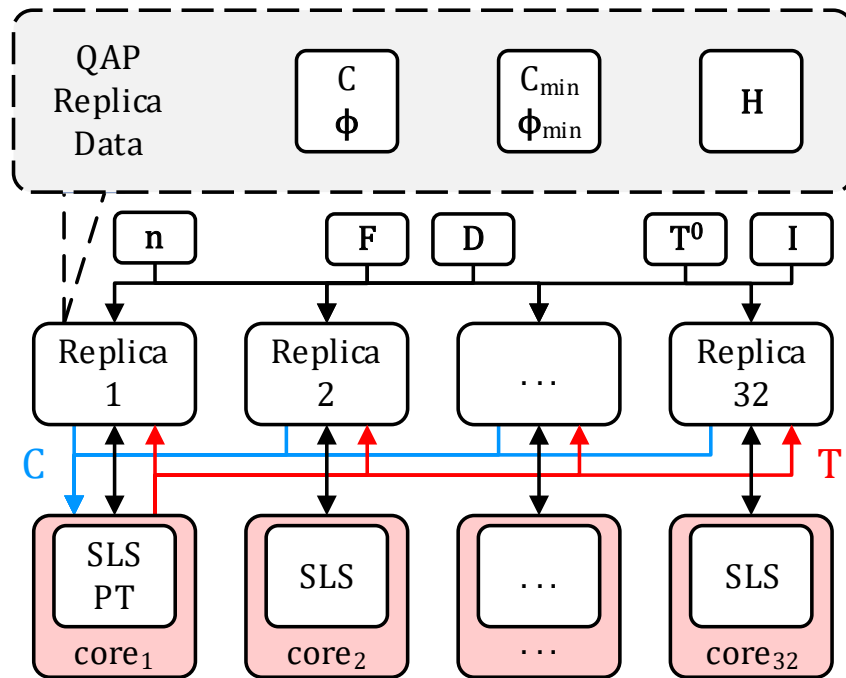
# Parallel Tempering



- Uses  $M$  unique search-states/replicas with unique temperatures,  $T = (T_k) \in R^M$
- Replicas perform Local Search at their assigned temperature
  - Can swap temperatures with other replicas with probability calculated via (17)
  - Allows for an escape mechanism from local minima

$$SAP = \min \left\{ 1, \exp \left( \left( \frac{1}{T_k} - \frac{1}{T_{k+1}} \right) (C_k - C_{k+1}) \right) \right\} \quad (17)$$

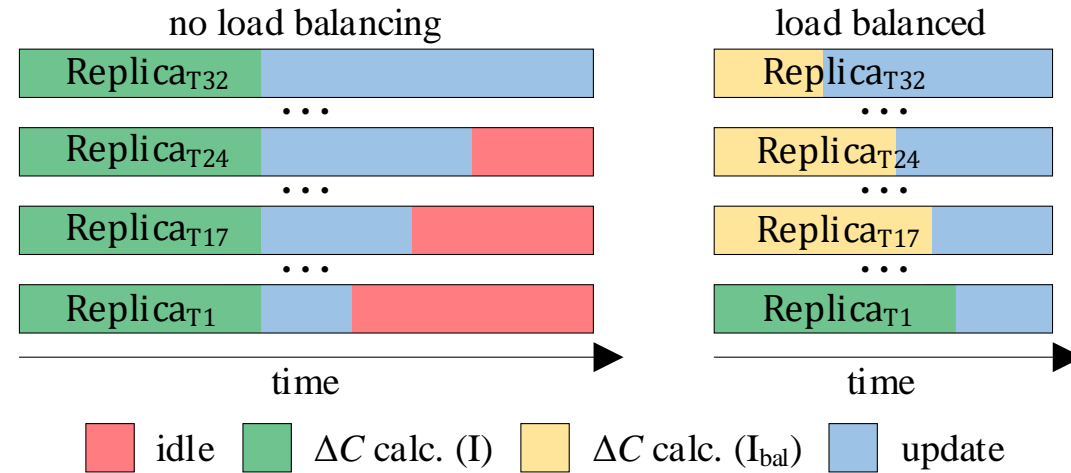
# Parallel Tempering Solver System



32 Replica PT Engine

- Experimental setup:
  - 64-Core AMD 3990X CPU with AVX2 support
  - Coded in C++ and compiled with GCC-9
- Float32 used to store  $F, D, H$ 
  - Allows for Fused-Multiply-Add instructions
- Each engine uses 32 replicas in combination with a PT controller
  - Each replica is run on a dedicated CPU core
  - 2 engines are run in parallel to utilize all 64 cores
  - Runs until a target cost is found or a 5-minute time-out limit is reached

# Load Balancing Threads



- Varying proposal acceptance rates of replicas at different temperatures causes thread imbalance when using Parallel Tempering
- Track average time-per-iteration at each temperature and scale iterations at each temperature to try to equalize thread run-times



# Outline

---

- Motivation and Background
  - Quadratic Assignment Problems
  - Local Search Algorithms
  
- Integer Boltzmann Machine
  
- Parallel Tempering Solver Design
  - Algorithm and SIMD
  - Load Balancing
  
- **Benchmark Results**
  
- Conclusion

# QAP Benchmark Results

TABLE I: QAP TtS Across Solvers (values in seconds)

Instance	ParEOTS	Int BM + PT	Speed-Up
dre42	0.7	0.10 ± 0.01	7.00×
dre56	5.6	0.99 ± 0.07	5.66×
dre72	26.0	5.23 ± 0.31	4.97×
Inst50	17.0	3.58 ± 0.27	4.75×
Inst60	67.0	2.00 ± 0.16	33.5×
Inst70	127.0	6.24 ± 0.98	20.35×
Inst80	116.0	19.36 ± 1.75	5.99×
sko81	22.4	0.13 ± 0.01	172.31×
<b>sko90</b>	<b>92.0</b>	<b>0.26 ± 0.01</b>	<b>353.85×</b>
sko100a	69.0	0.27 ± 0.01	255.56×
sko100b	45.0	0.19 ± 0.01	236.84×
sko100c	56.0	0.29 ± 0.02	193.1×
sko100d	37.0	0.38 ± 0.02	97.37×
sko100e	47.0	0.21 ± 0.01	223.81×
sko100f	57.0	0.24 ± 0.01	237.5×
tai60b	4.6	0.14 ± 0.01	32.86×
tai80b	53.0	0.43 ± 0.01	123.26×
tai100b	71.0	0.36 ± 0.02	197.22×
wil100	97.0	0.46 ± 0.02	210.87×
<b>Geometric Mean Speedup</b>			<b>57.09×</b>

- Avg. Time-to-Solution (TtS) across 100 independent, randomly seeded, runs with 99% confidence interval for 19 difficult QAP instances
  - Instances from [4,5,6]
  
- ParEOTS [7]: Parallel Extremal Optimization + Tabu Search
  - Best performing published QAP solver
  - 128 cores on CPU cluster 8 x 16-core AMD 6376 CPU
  - TtS reported as average across 10 runs with 5-minute time-out
  
- 57x faster than ParEOTS on average
  - Max speed-up of 353x

# Conclusion

---

- Extended BMs to support for integer variables with one-hot constraints
  - Remove need for penalty tuning on user side
- Added support for weight-generation using submatrices for problem topologies such as QAP
  - Larger problem support due to lower memory requirements
- Integer BM based CPU solution displays performance that is orders of magnitude faster than competing solvers across difficult QAP instances

## Future Work

- Accelerate search relative to high-performance multi-core CPU system through dedicated hardware

# Acknowledgements

---

The authors would like to thank Fujitsu Ltd. and Fujitsu Consulting (Canada) Inc. for providing financial support and technical expertise on this research and to CMC Microsystems for access to CAD tools used in this research

# Thank You!

---

Looking forward to your questions at the live session

# References

---

- [1] Lawler, E.L.: The quadratic assignment problem. *Manage. Sci.* 9(4), 586–599 (1963).
- [2] M. Bagherbeik, P. Ashtari, S. F. Mousavi, K. Kanda, H. Tamura, and A. Sheikholeslami, “A permutational boltzmann machine with parallel tempering for solving combinatorial optimization problems,” in *International Conference on Parallel Problem Solving from Nature*. Springer, 2020, pp. 317–331
- [3] M. Bagherbeik and A. Sheikholeslami, “Caching and vectorization schemes to accelerate local search algorithms for assignment problems,” in *IEEE Congress on Evolutionary Computation (to appear)*, 2021.
- [4] Burkard, R.E., Karisch, S.E., Rendl, F.: QAPLIP-a quadratic assignment problem library. *J. Global Optim.* 10(4), 391–403 (1997).
- [5] Palubeckis, G.: An algorithm for construction of test cases for the quadratic assignment problem. *Informatika Lith. Acad. Sci.* 11, 281–296 (2000)
- [6] Drezner, Z., Hahn, P.M., Taillard, E.D.: Recent advances for the quadratic assignment problem with special emphasis on instances that are difficult for metaheuristic methods. *Ann. Oper. Res.* 139(1), 65–94 (2005).
- [7] Munera, Danny, Daniel Diaz, and Salvador Abreu. "Hybridization as cooperative parallelism for the quadratic assignment problem." *International Workshop on Hybrid Metaheuristics*. Springer, Cham, 2016.