



uDeviceX



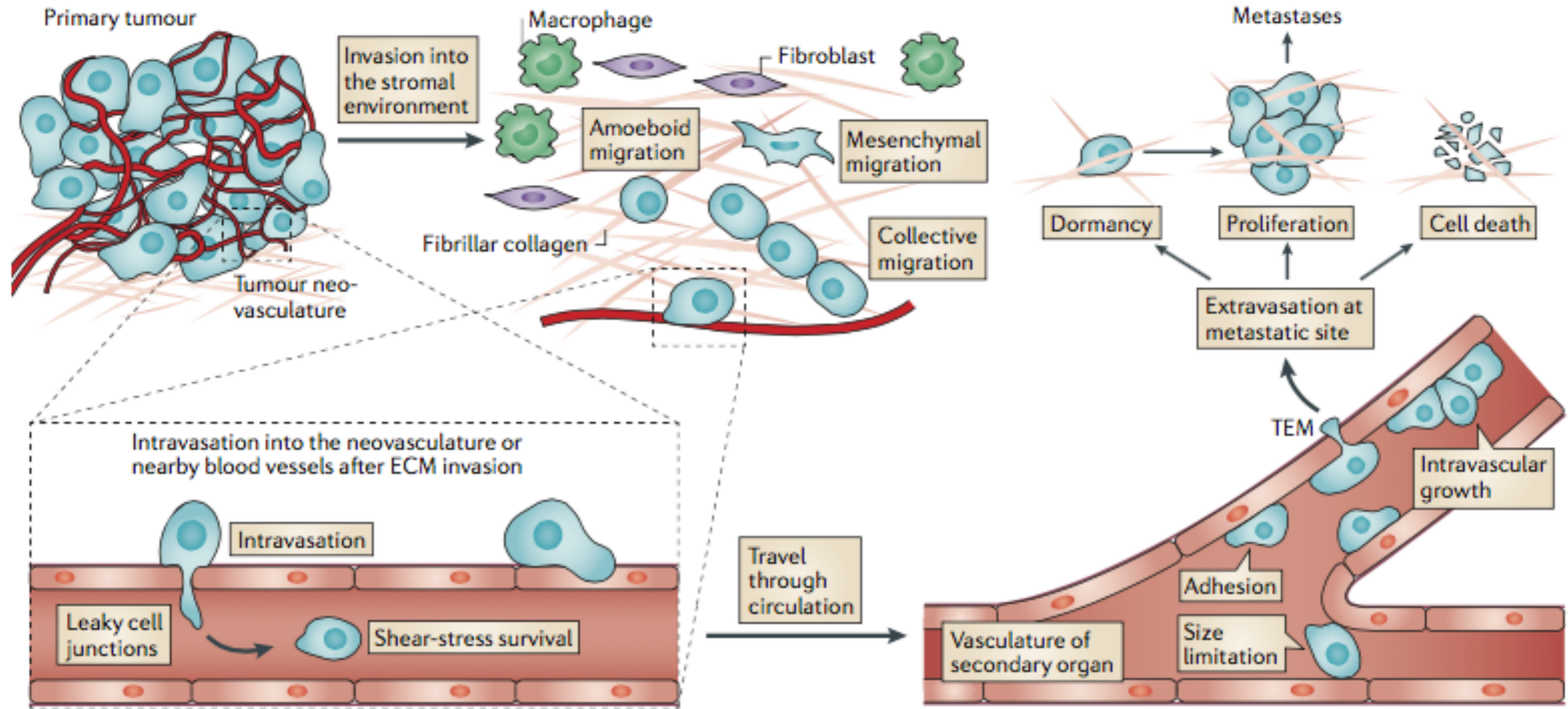
18'688 K20Xs

# RUNNING AFTER A TUMOR CELL

Diego Rossinelli, ETH Zurich

The In-Silico Lab-on-a-Chip: Petascale and High-Throughput Simulations of Microfluidics at Cell Resolution  
Rossinelli, Tang, Lykov, Alexeev, Bernaschi, Hadjidoukas, Bisson, Joubert, Conti, Karniadakis, Fatica, Pivkin, Koumoutsakos  
ACM Gordon Bell finalist 2015

# METASTASIS AND CTCs



NATURE REVIEWS | **CANCER**

VOLUME 13 | DECEMBER 2013 | 859

- 8 million cancer deaths every year
- 90% attributed to metastasis

CTCs in the blood  
-> higher risks of metastasis

How to detect the presence of CTCs in a blood sample?

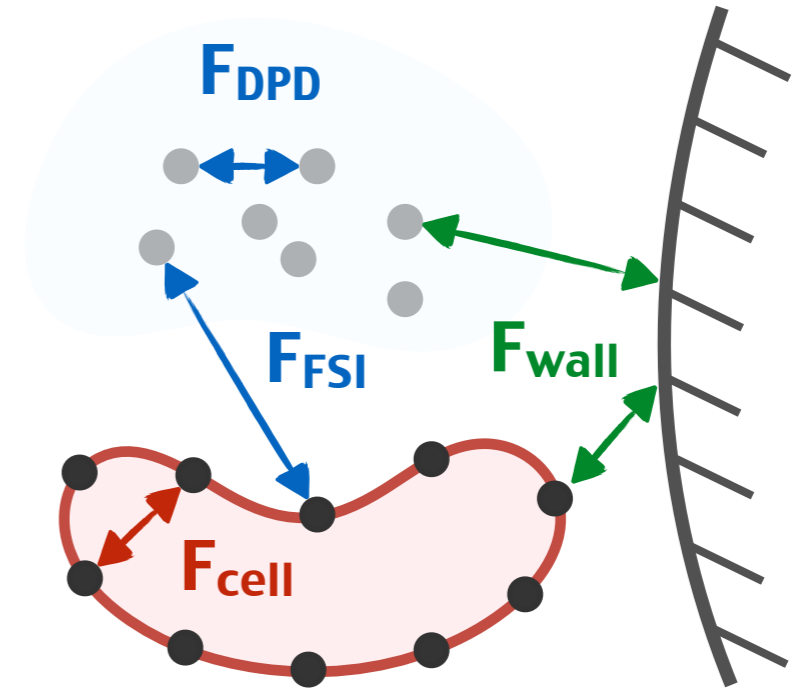
# THE CTC-ICHIP

Circulating Tumor Cells  
(CTCs)

# NUMERICAL METHOD

# DISSIPATIVE PARTICLE DYNAMICS

$$\begin{aligned}
 \mathbf{F}_i = & \sum_{n=1, n \neq i}^N \mathbf{F}_{i,n}^{C,DPD} + \mathbf{F}_{i,n}^{D,DPD} + \mathbf{F}_{i,n}^{R,DPD} \\
 & + \sum_{k=1, k \neq i}^K \mathbf{F}_{i,k}^{C,FSI} + \mathbf{F}_{i,k}^{D,FSI} + \mathbf{F}_{i,k}^{R,FSI} \\
 & + \sum_{m=1, m \neq i}^M \mathbf{F}_{i,m}^{C,wall} + \mathbf{F}_{i,m}^{D,wall} + \mathbf{F}_{i,m}^{R,wall}
 \end{aligned}$$



$$\mathbf{F}_{ij}^C = \begin{cases} a_{ij}(1 - r_{ij})\mathbf{e}_{ij}, & \text{if } r_{ij} < 1 \\ 0, & \text{if } r_{ij} \geq 1 \end{cases}$$

$$\mathbf{F}_{ij}^D = -\gamma w^D(r_{ij})(\mathbf{e}_{ij} \cdot \mathbf{v}_{ij})\mathbf{e}_{ij}$$

$$\mathbf{F}_{ij}^R = \sigma w^R(r_{ij})\theta_{ij}\mathbf{e}_{ij}$$

$$\mathbf{r}_{ij} = \mathbf{r}_i - \mathbf{r}_j$$

$$r_{ij} = \|\mathbf{r}_{ij}\|$$

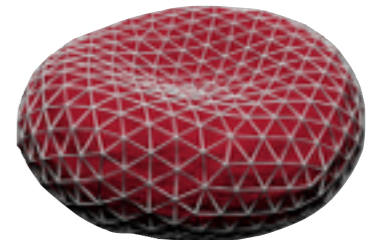
$$\mathbf{e}_{ij} = \mathbf{r}_{ij} / \|\mathbf{r}_{ij}\|$$

$$\mathbf{v}_{ij} = \mathbf{v}_i - \mathbf{v}_j$$

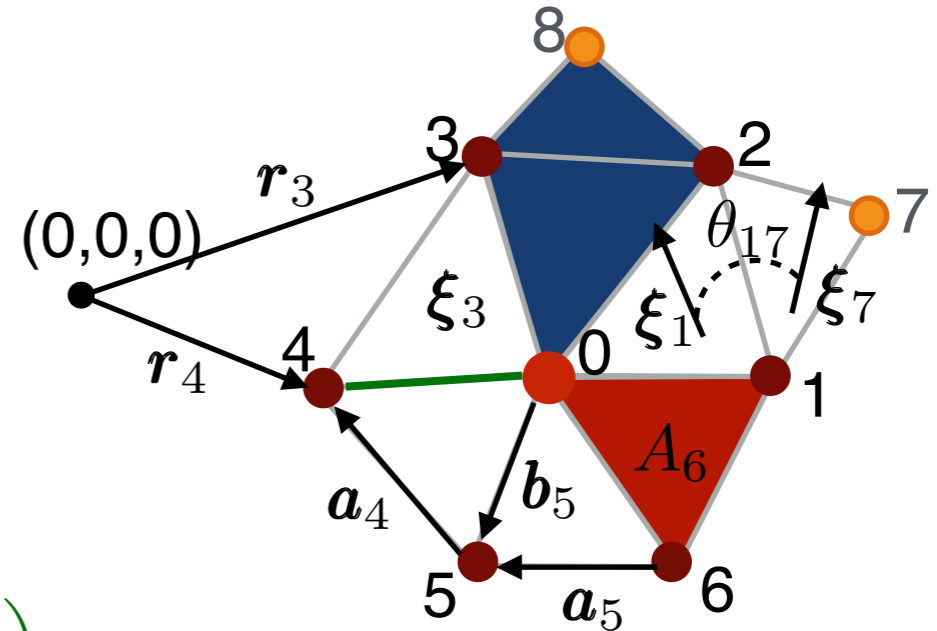
$$w^D(r) = (w^R(r))^2$$

$$\sigma^2 = 2\gamma k_B T$$

# FORCES IN THE SYSTEM: RBC



$$\begin{aligned}
 \mathbf{F}^{\text{cell}} &= \sum_{n=1}^N \mathbf{F}_{0,n-1,n,n+1}^{\text{dihedral},1} + \mathbf{F}_{0,n,N+n,n+1}^{\text{dihedral},2} + \mathbf{F}_{0,n,n+1}^{\text{triangle}} + \mathbf{F}_{0,n}^{\text{bond}} \\
 &= \sum_{n=1}^N \beta_{n,n+1}^b \left[ \frac{\boldsymbol{\xi}_n \times \mathbf{a}_{n+1} + \boldsymbol{\xi}_{n+1} \times \mathbf{a}_n}{\xi_n \xi_{n+1}} - \cos \theta_{n,n+1} \left( \frac{\boldsymbol{\xi}_n \times \mathbf{a}_n}{\xi_n^2} + \frac{\boldsymbol{\xi}_{n+1} \times \mathbf{a}_{n+1}}{\xi_{n+1}^2} \right) \right] \\
 &+ \beta_{n,N+n}^b \left[ \frac{\boldsymbol{\xi}_{N+n} \times \mathbf{a}_n}{\xi_n \xi_{N+n}} - \cos \theta_{n,N+n} \frac{\boldsymbol{\xi}_n \times \mathbf{a}_n}{\xi_n^2} \right] \\
 &+ \left( \frac{qC_q}{A_n^{q+1}} - k_a \frac{A - A_0^{\text{tot}}}{A_0^{\text{tot}}} \right) \frac{\boldsymbol{\xi}_n \times \mathbf{a}_n}{4A_n} \\
 &- \frac{k_v}{18} \frac{V - V_0^{\text{tot}}}{V_0^{\text{tot}}} (\boldsymbol{\xi}_n + (\mathbf{r}_0 + \mathbf{r}_n + \mathbf{r}_{n+1}) \times \mathbf{a}_n) \\
 &- \frac{k_B T}{p} \left( \frac{1}{4(1 - b_n/l_m)^2} - \frac{1}{4} + \frac{b_n}{l_m} \right) \frac{\mathbf{b}_n}{b_n} \\
 &+ \sqrt{2k_B T} \left( \sqrt{2\gamma^T} d\overline{\mathbf{W}}_{ij}^S + \sqrt{3\gamma^C - \gamma^T} \frac{\text{tr}[d\mathbf{W}_{ij}]}{3} \mathbf{I} \right)
 \end{aligned}$$



$$x_0 = l_0/l_m$$

$$A_0 = \sqrt{3}l_0^2/4$$

$$C_q = \frac{\sqrt{3}A_0^{q+1}k_B T(4x_0^2 - 9x_0 + 6)}{4pql_m(1 - x_0)^2}$$

$$\beta_{ij}^b = k_b \frac{\sin \theta_{ij} \cos \theta_0 - \cos \theta_{ij} \sin \theta_0}{\sqrt{1 - \cos^2 \theta_{ij}}}$$

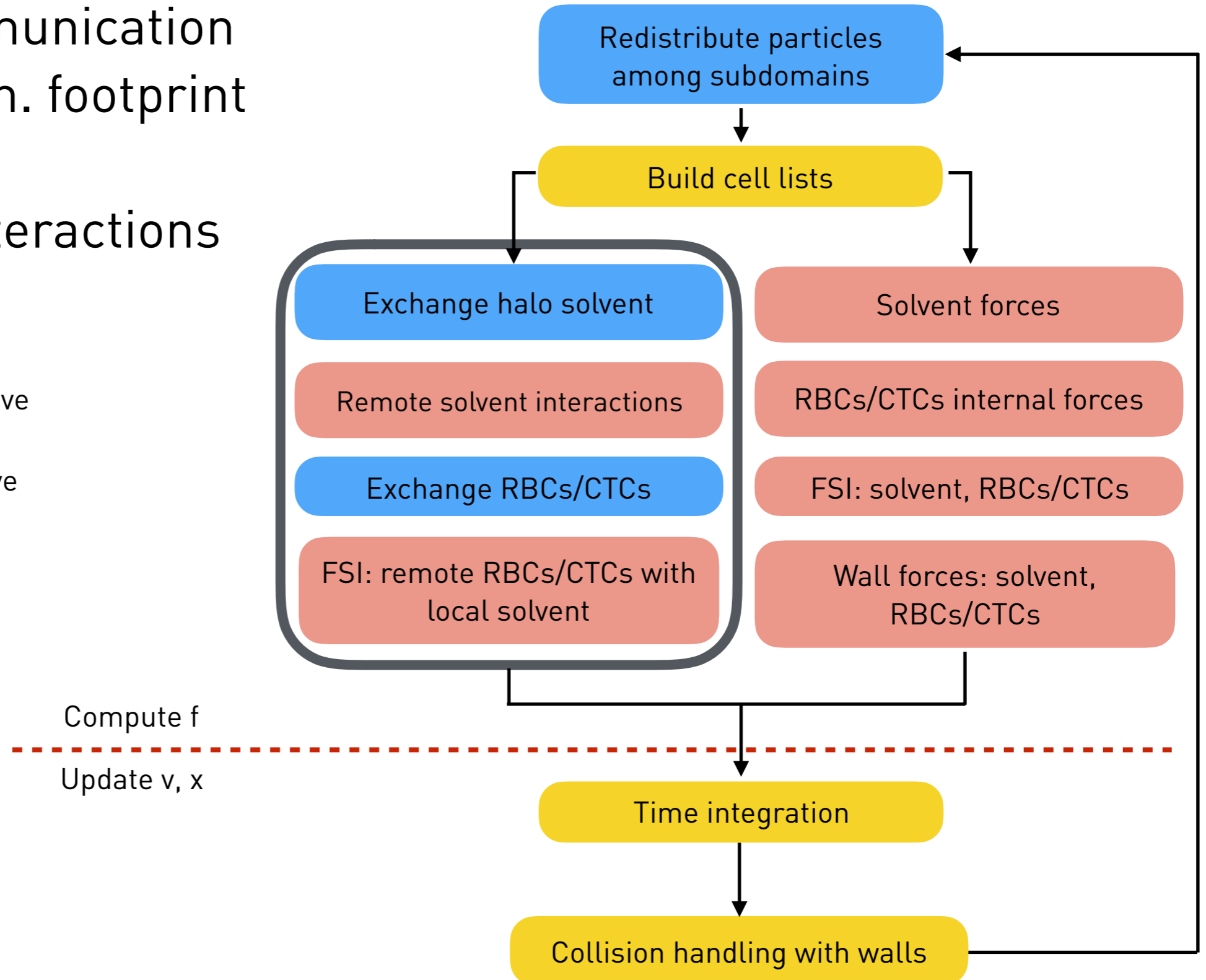
$$d\overline{\mathbf{W}}_{ij}^S = d\mathbf{W}_{ij}^S - \text{tr}[d\mathbf{W}_{ij}^S] \mathbf{1}/3$$

# ONE TIME STEP

Non-trivial communication

- 15-30MB comm. footprint
- 180 messages
- 5-40 million interactions

- Compute intensive
- Memory intensive
- Communication
- Remote forces



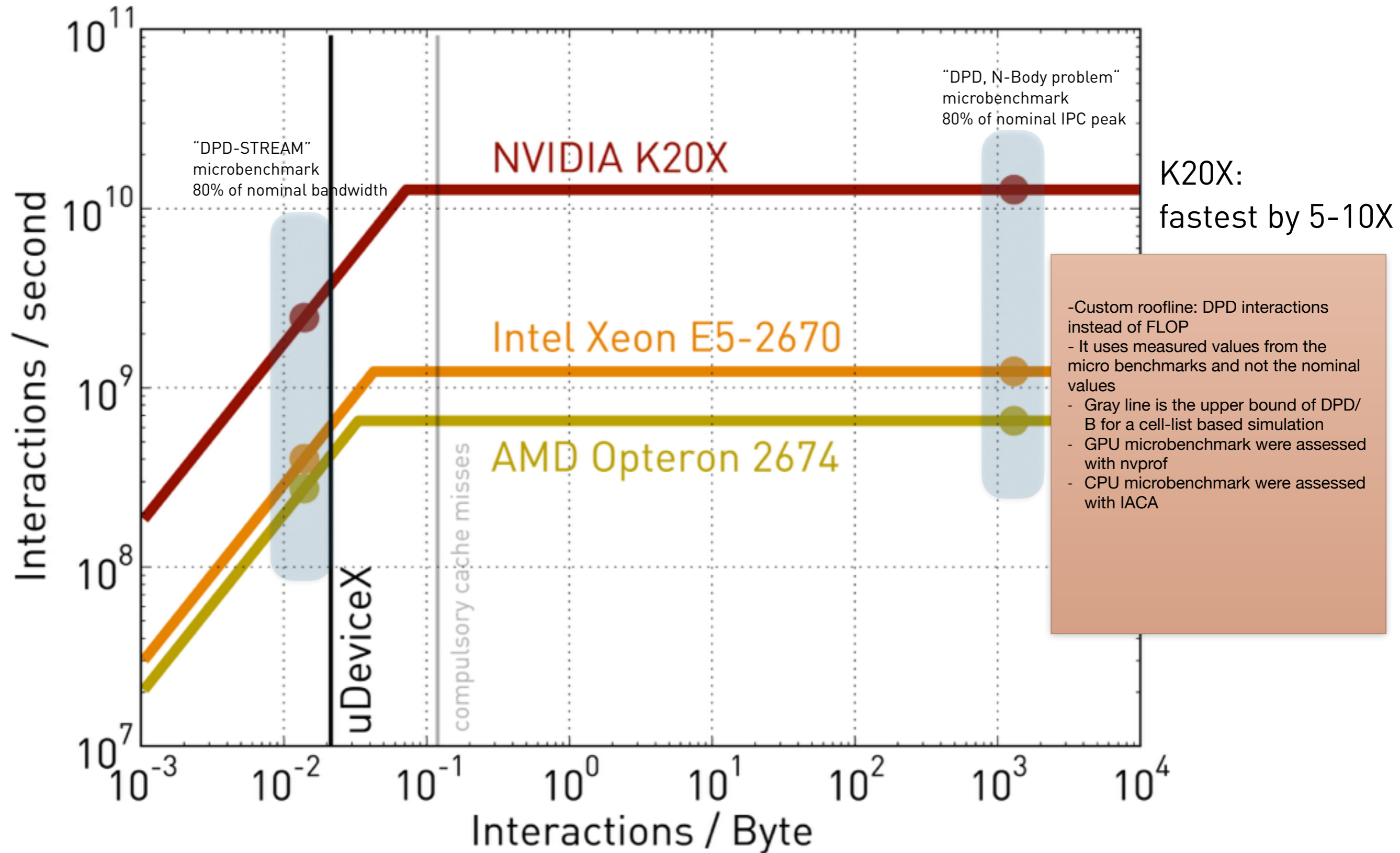


# HPC

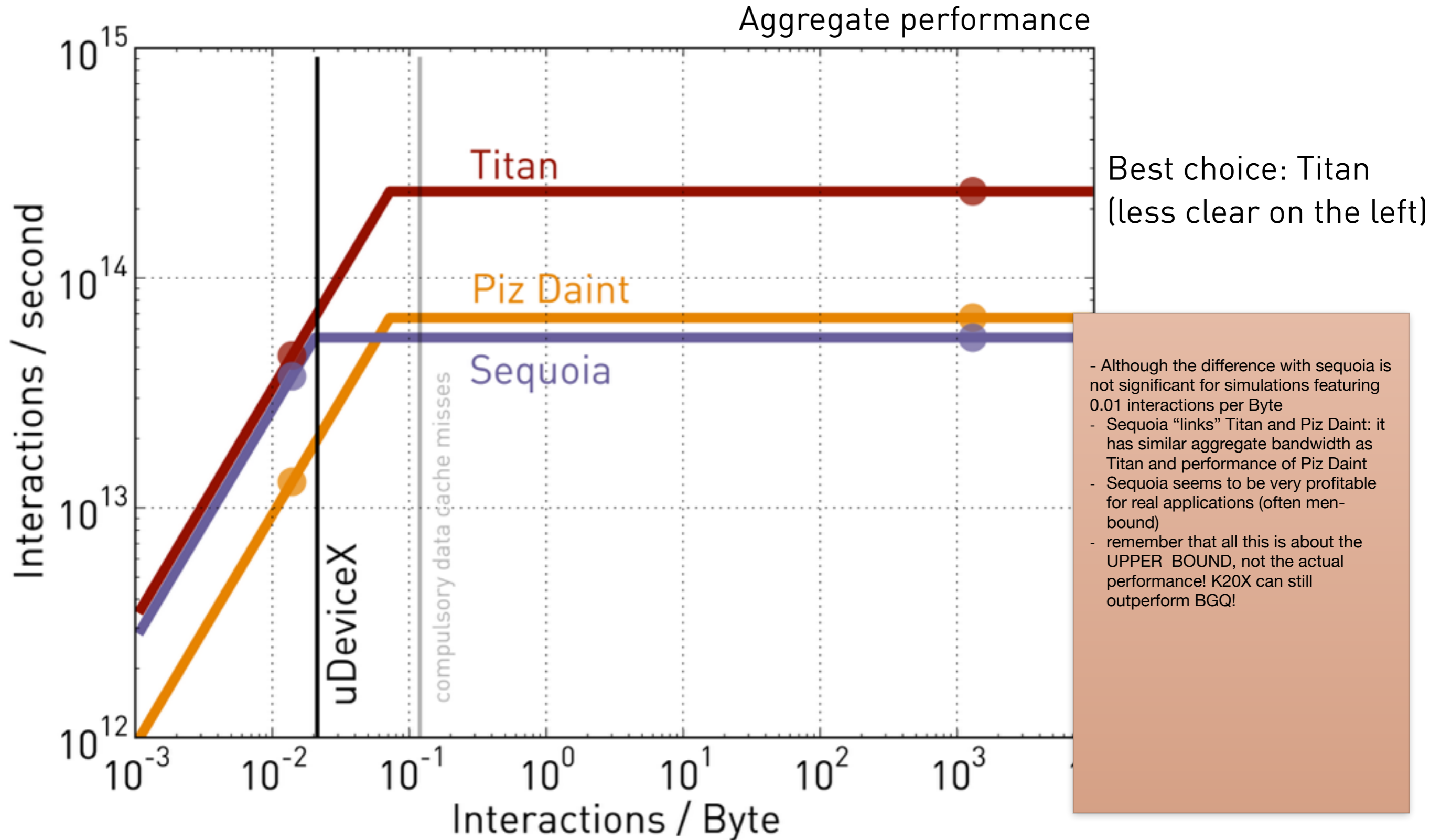
Minimize #interactions

Maximize interactions/s  $\rightarrow$  Maximize IPC

# DPD: HARDWARE OF CHOICE



# DPD: SUPERCOMPUTER OF CHOICE



# SUPERCOMPUTING CHALLENGES

see A28 and A29 of Patterson

$\times \#instr.$

$1/IPC$

$TTS = 1/freq \times$

## DPD

## CHALLENGE

Rapidly changing neighbours

Expensive forces

No gain from Verlet lists

High instruction count

IOP-based random forces

Irregular computation  
Irregular access patterns  
Irregular inter-rank  
messages

Low IPC on the K20X

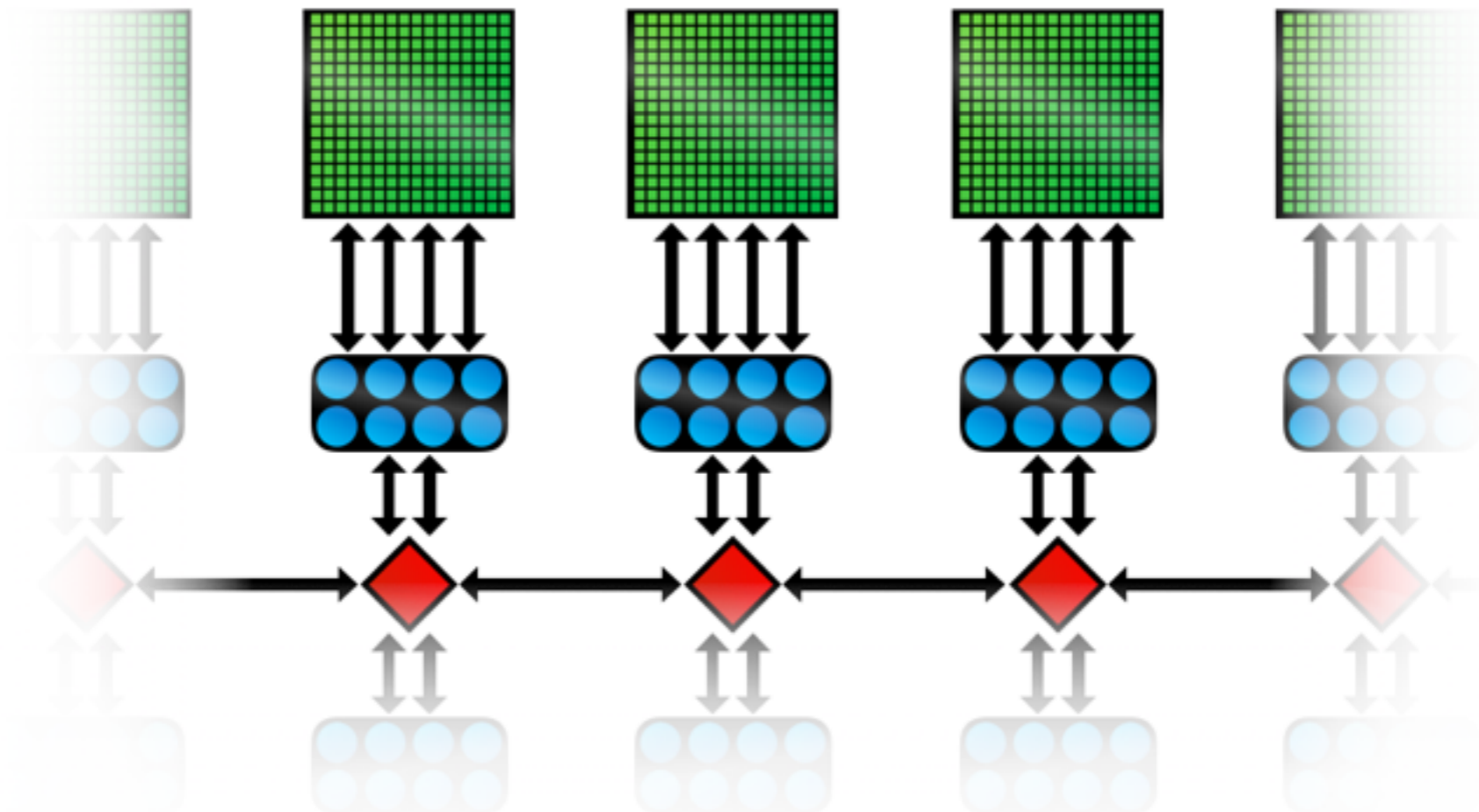
Warp divergence  
Uncoalesced access  
Penalised network performance  
Poor C/T overlap

➔ **Latency-bound performance**

➔ **Low IPC**

# HOW TO MAXIMISE IPC

---

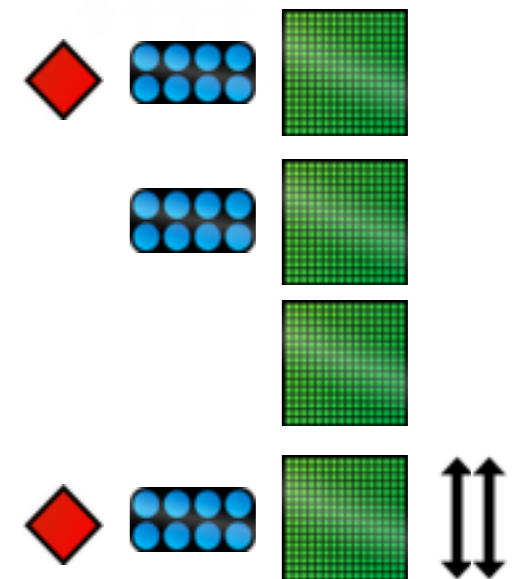


➔ Relax irregularities

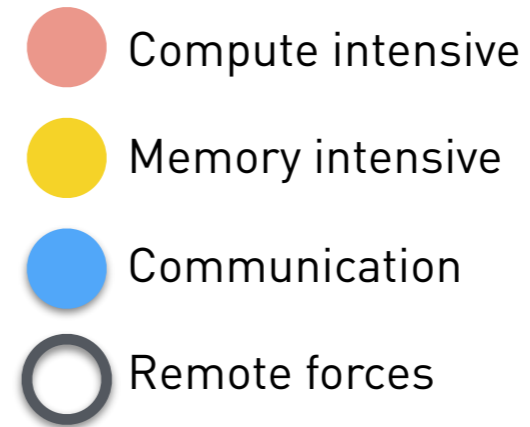
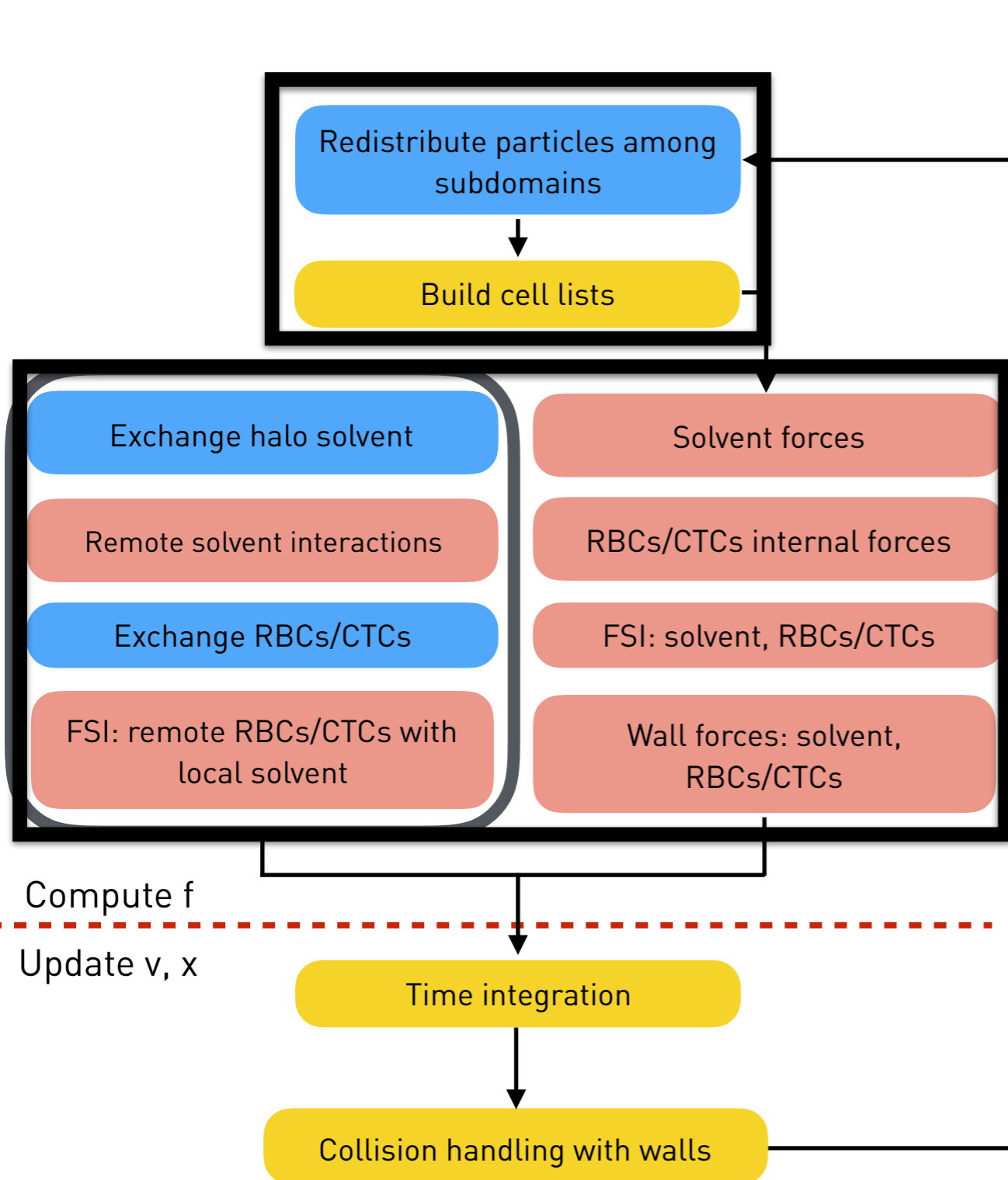
➔ Balance the workload

➔ Prefer high-throughput instructions

➔ Exploit the parallelism in the system



# NETWORK LEVEL

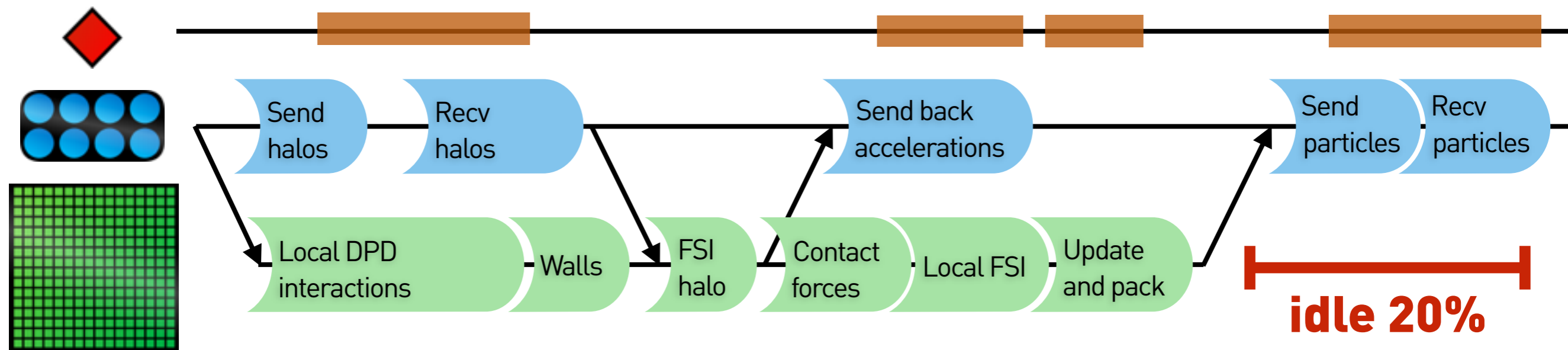


- Message sizes guessed a-priori
- Exceptions: secondary messages
- Adapted over recent history

- ➔ Non-blocking MPI calls
- ➔ C/T Overlap

# NODE LEVEL

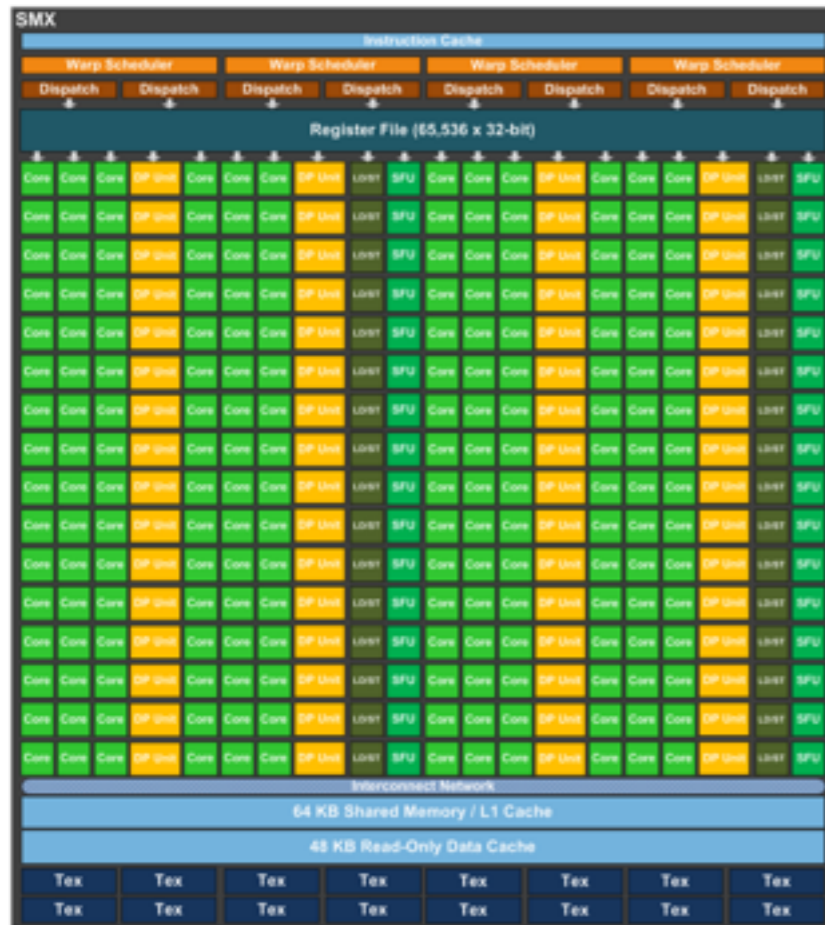
- Fully asynchronous CPU-GPU workflow
- Non-trivial dependency graph with 3 CUDA streams



Solution:

- Node oversubscription with multiple MPI tasks
- ➔ **GPU utilisation goes to over 95%**
- ➔ **Opportunities for load-balancing**

# K20X FOR DPD SIMULATIONS



Throughput of Native Arithmetic Instructions.  
(Number of Operations per Clock Cycle per Multiprocessor)

	Compute Capability					
	2.0	2.1	3.0, 3.2	3.5, 3.7	5.0, 5.2	5.3
16-bit floating-point add, multiply, multiply-add	N/A	N/A	N/A	N/A	N/A	256
32-bit floating-point add, multiply, multiply-add	<b>SP: best operation</b>			192	128	128
64-bit floating-point add, multiply, multiply-add	<b>DP: 3x penalty</b>			$64^2$	4	4
32-bit floating-point reciprocal, reciprocal square root, base-2 logarithm ( <code>__log2f</code> ), base 2 exponential ( <code>exp2f</code> ), sine ( <code>__sinf</code> ), cosine ( <code>__cosf</code> )	4	8	32	32	32	32
32-bit integer add, extended-precision add, subtract, extended-precision subtract	32	<b>OK</b>	160	160	128	128
32-bit integer multiply, multiply-add, extended-precision multiply-add	<b>Integer multiplication: 6x penalty</b>			32	Multiple instructions	Multiple instructions
24-bit integer multiply ( <code>__[u]mu124</code> )	Multiple instructions	Multiple instructions	Multiple instructions	Multiple instructions	Multiple instructions	Multiple instructions
32-bit integer shift	16	16	32	$64^3$	64	64
compare, minimum, maximum	32	<b>OK</b>	160	160	64	64
32-bit integer bit reverse, bit field extract/insert	16	16	32	32	64	64
32-bit bitwise AND, OR, XOR	32	<b>OK</b>	160	160	128	128
count of leading zeros, most significant non-sign bit	16	16	32	32	Multiple instructions	Multiple instructions

➔ Maximize IPC

➔ **Maximize GPU throughput**



# RANDOM NUMBER GENERATOR

---

- FMA-based RNG
- Passes BigCrush TestU01
- At least 18 rounds

	SARU	OURS
FP32	30	<b>64</b>
non-FP32	81	16
TOT	111	<b>80</b>

```
function MEAN0VAR1(i,j,k)
    // Low-discrepancy number
    u ← MIN(i,j)
    v ← MAX(i,j)
    pij ← MOD(u × G + v × S,1)
    y ← k - pij

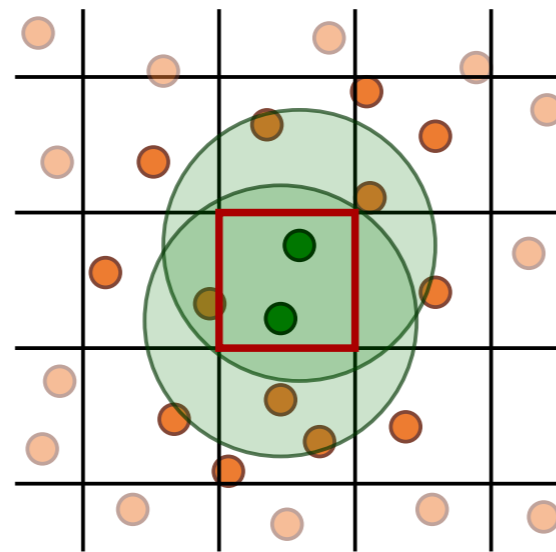
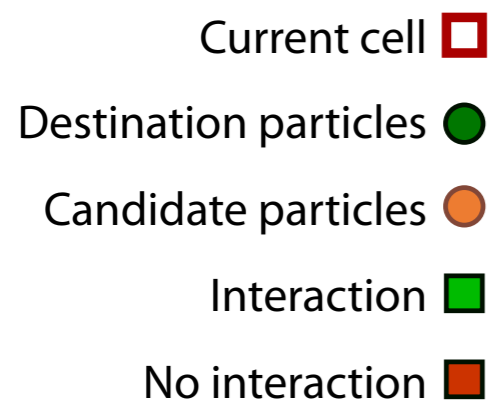
    // Pass through Logistic map
    for N rounds do
        y ← 4y(1 - y)
    end for


    // Normalize
    z ← NORMALIZE(y)
    return z
end function
```

# GPU LEVEL

Verlet lists:

- **Produced and consumed on-chip**
- Front-end: Verlet lists production
- Backend: DPD Interactions

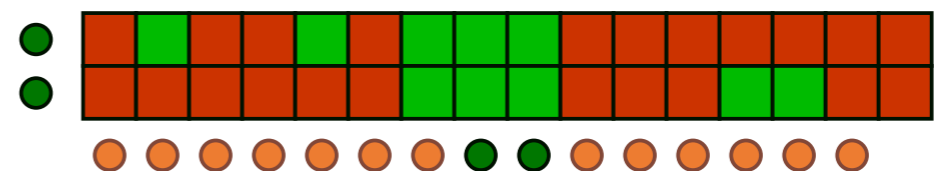


Cell	NW	N	NE	W	Self	E	SW	S	SE
# of 	1	1	3	2	2	1	1	3	1

Prefix sum

0	1	2	5	7	9	10	11	14	15
---	---	---	---	---	---	----	----	----	----

Map to 32 threads



Workload of a cell-list dynamically mapped to a warp

- Latencies of the front-end hidden by the backend
- Warp dynamically configured (1x32, 2x16, 4x8)
- **SIMT far more flexible than SIMD here!**

# uDeviceX

<http://udevicex.github.io/uDeviceX/>

# THE IN SILICO CTC-ICHIP



In-Vitro

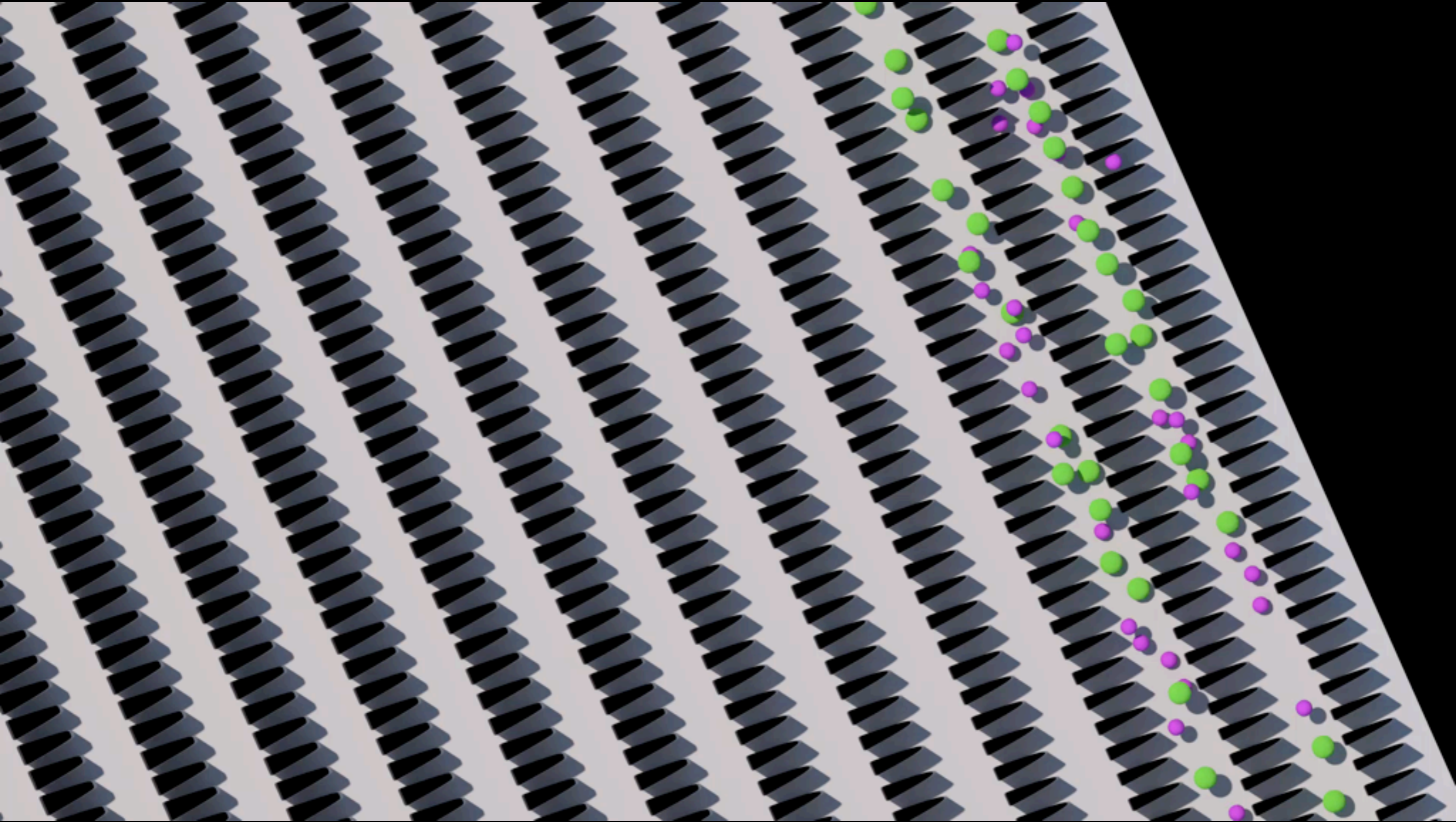
- **TIME TO SOLUTION:** >45X over State of the Art
- **SCALING:** >98% Weak and >87% Strong
- **PERFORMANCE:** max: 66% of nominal peak (avg. 34%)



# COMPLEXITY: 1 to 1 with microfluidic Devices at subcell resolution

- 1.0 E+13 DPD particles
- 1.0 E+08 Time Steps
- Timescale: 10 seconds
- 0.3 ml Blood - 1.4 Billion RBCs
- $\mu$ Fluidic Chips up to 50 mm<sup>3</sup>

# FUNNEL RATCHETS DEVICE



# ACKNOWLEDGMENTS

---

**ETH**

Eidgenössische Technische Hochschule Zürich  
Swiss Federal Institute of Technology Zurich



**BROWN**

Università  
della  
Svizzera  
italiana



Consiglio  
Nazionale delle  
Ricerche



**CSCS**

Centro Svizzero di Calcolo Scientifico  
Swiss National Supercomputing Centre



National Institutes  
of Health

**FNSNF**

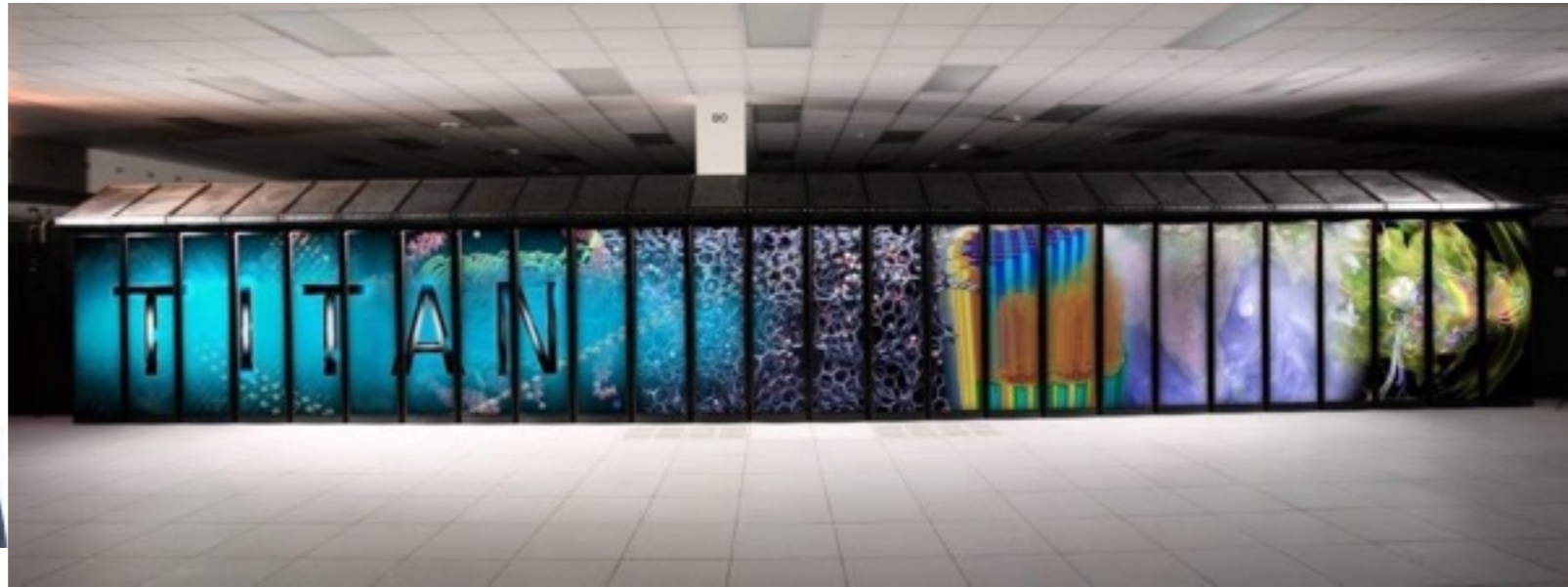
FONDS NATIONAL SUISSE  
SCHWEIZERISCHER NATIONALFONDS  
FONDO NAZIONALE SVIZZERO  
SWISS NATIONAL SCIENCE FOUNDATION



European Research Council

Established by the European Commission





THE END

Thank you!

# DISSIPATIVE PARTICLE DYNAMICS

- Models the solvent
- Short-range interactions
- Fluctuation-dissipation theorem
- Correct hydrodynamics for  $Re < 10$

$$\mathbf{F}_i = \sum_{j \neq i}^N (\mathbf{F}_{ij}^C + \mathbf{F}_{ij}^D + \mathbf{F}_{ij}^R), \quad i = 1, \dots, N$$

**Conservative:**

$$\mathbf{F}_{ij}^C = \sigma w(r_{ij}) \frac{\mathbf{r}_{ij}}{r_{ij}}$$

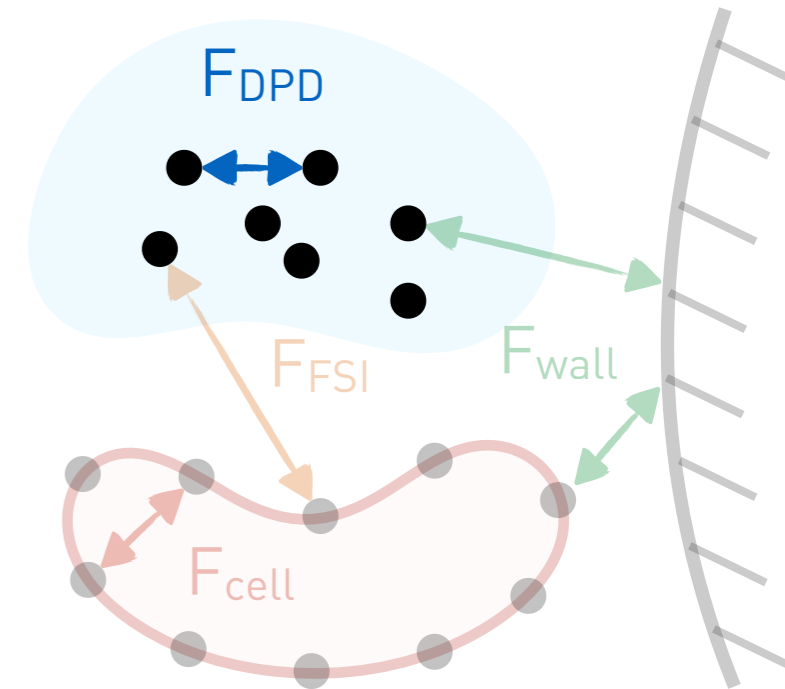
**Random:**

$$\mathbf{F}_{ij}^R = \sigma w(r_{ij}) \theta \frac{\mathbf{r}_{ij}}{r_{ij}}$$

Kernel

**Dissipative:**

$$\mathbf{F}_{ij}^D = -\gamma (w(r_{ij}))^2 (\mathbf{r}_{ij}, \mathbf{u}_{ij}) \frac{\mathbf{r}_{ij}}{r_{ij}^2}$$



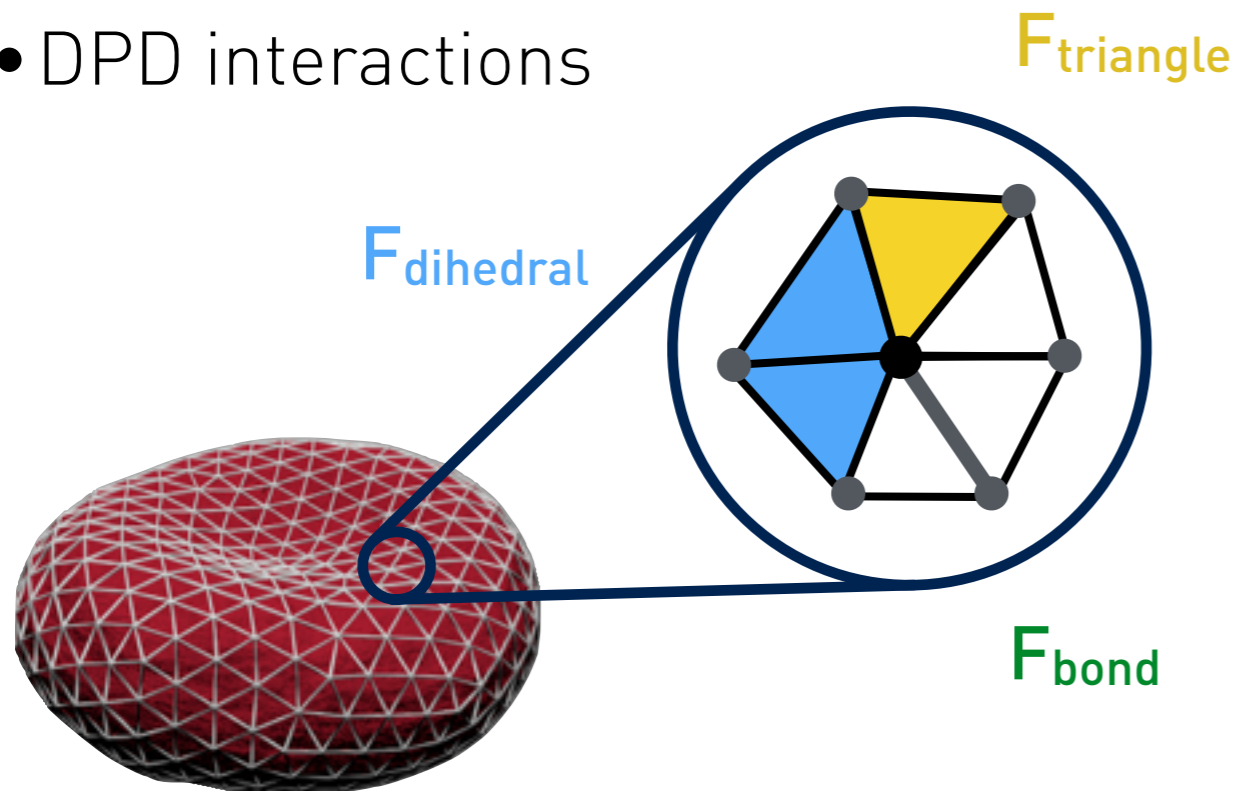
# CELLS AND FLUID-STRUCTURE INTERACTIONS

## Cells:

- Deformable membrane model
- Discretised as triangle mesh
- Membrane forces are stiff

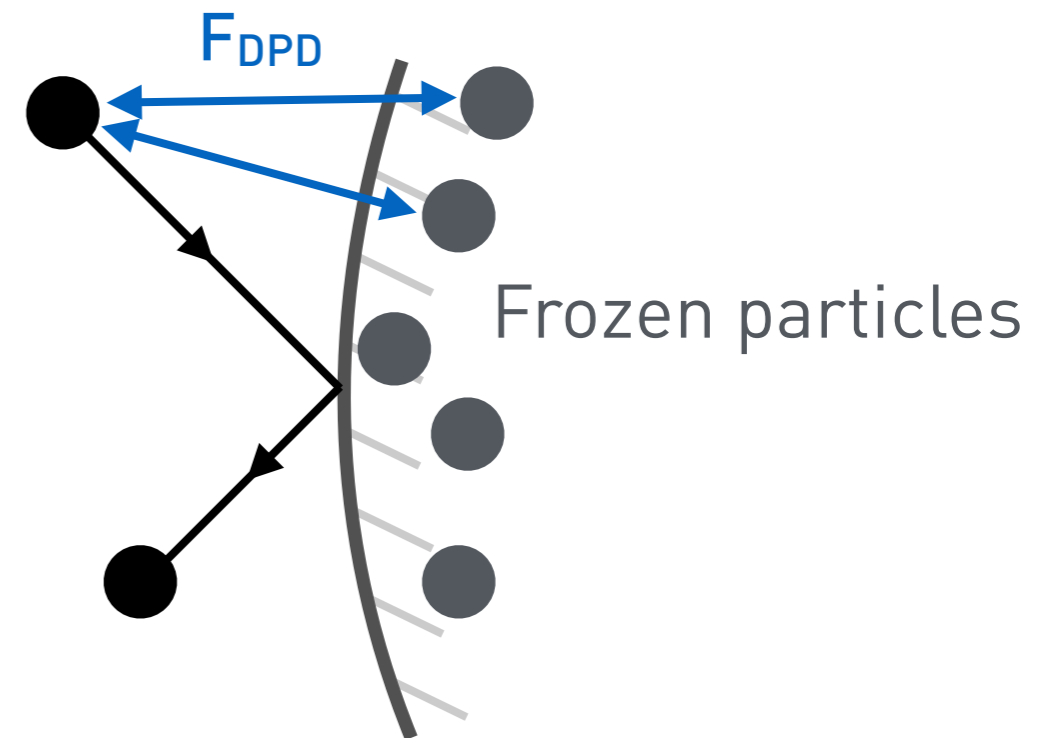
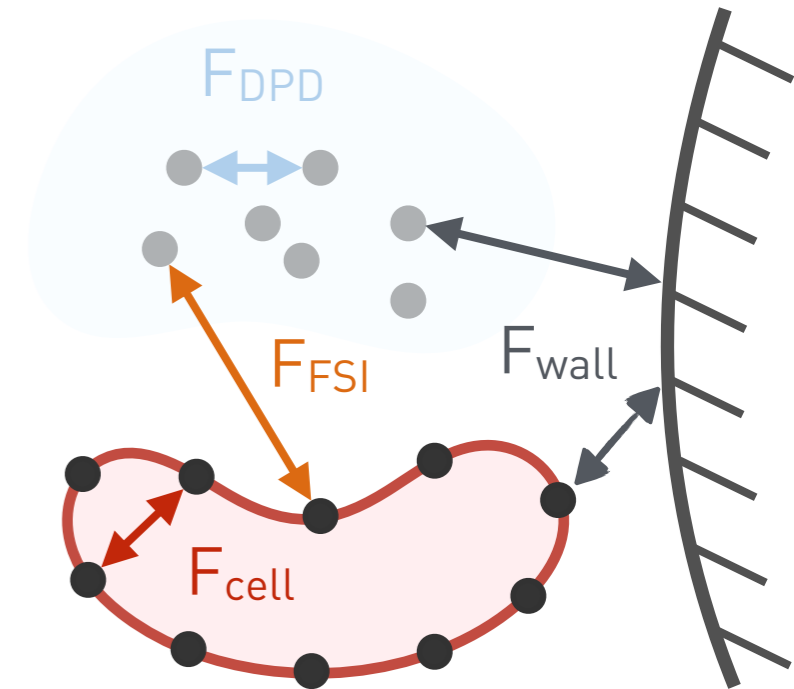
## FSI forces:

- DPD interactions

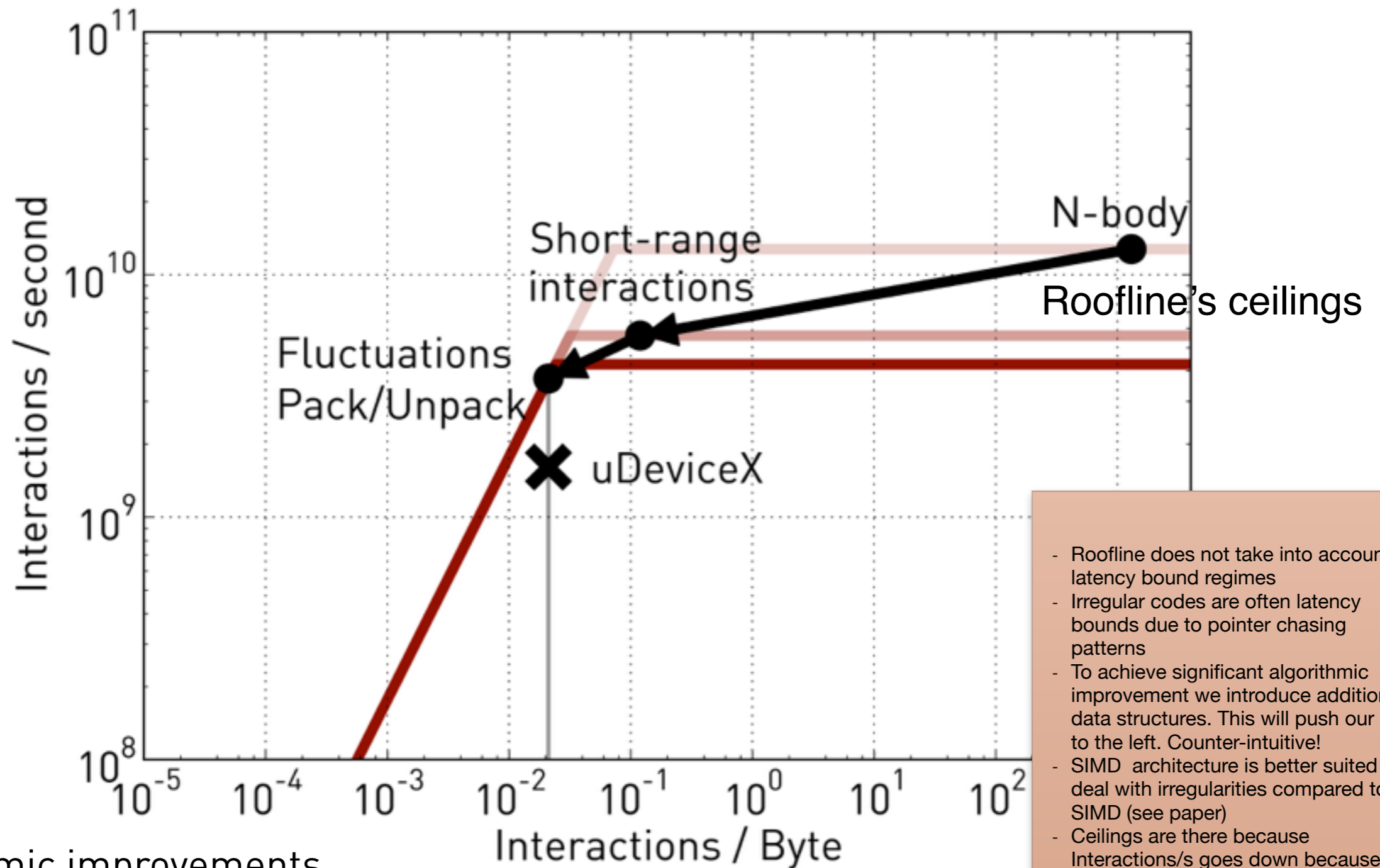


## Walls:

- Implicit description
- Frozen particles
- Bounce-back



# THE DOWNSIDE OF ALGORITHMIC IMPROVEMENTS

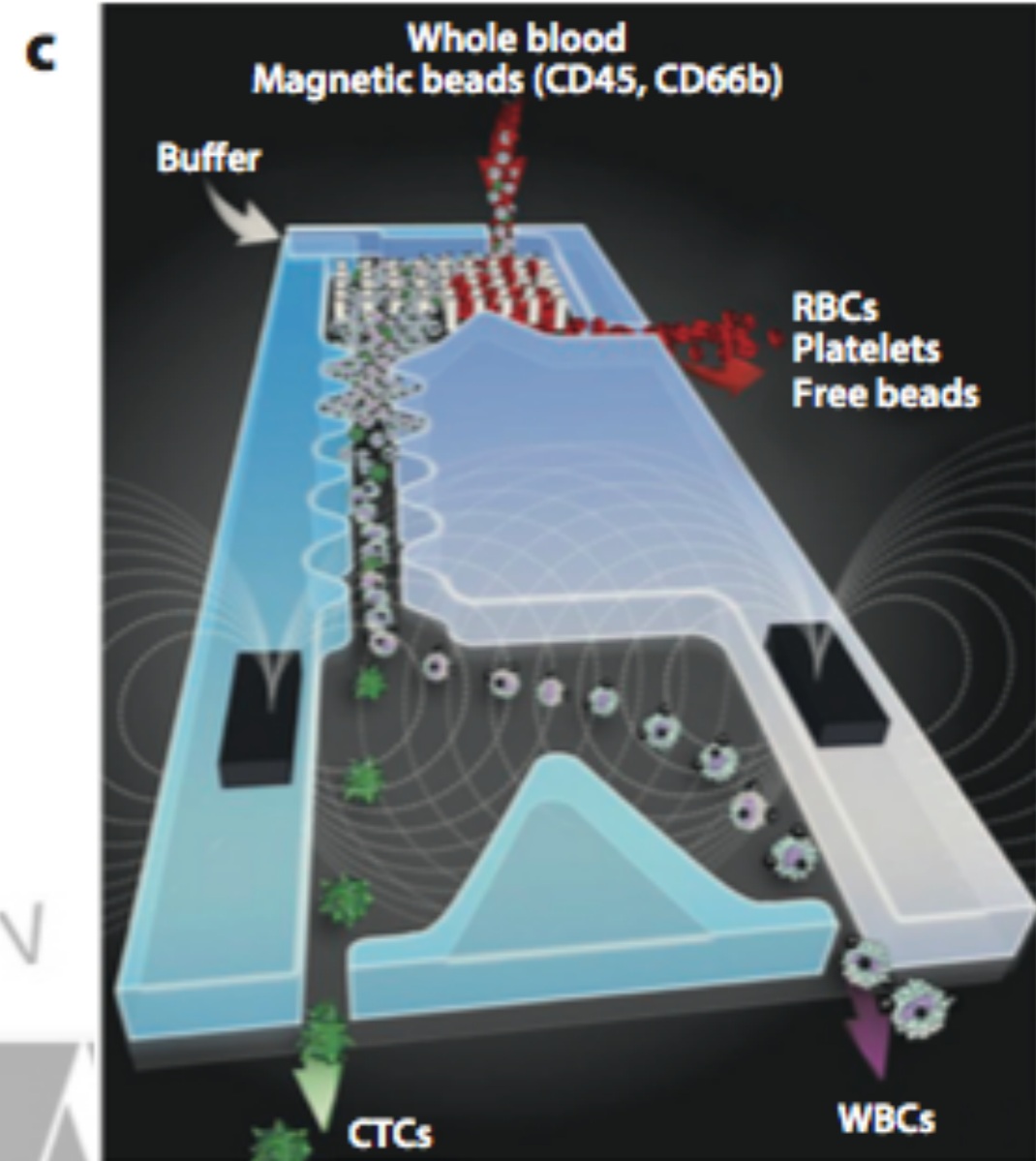
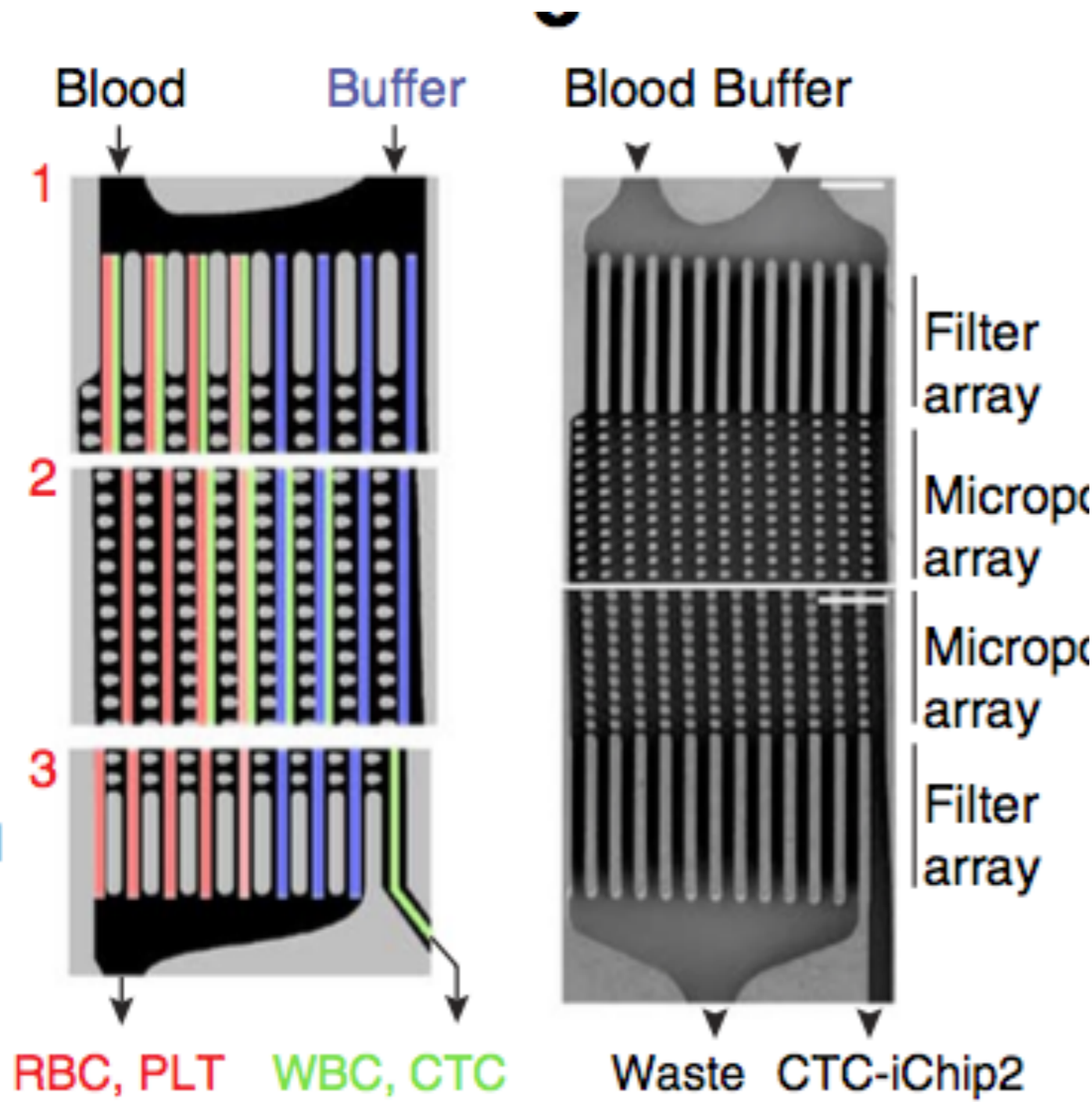


## Algorithmic improvements

- Lead to irregular codes
- Often enters in a latency-bound regime
- OIs shift to the left

- Roofline does not take into account latency bound regimes
- Irregular codes are often latency bounds due to pointer chasing patterns
- To achieve significant algorithmic improvement we introduce additional data structures. This will push our OI to the left. Counter-intuitive!
- SIMD architecture is better suited to deal with irregularities compared to SIMD (see paper)
- Ceilings are there because Interactions/s goes down because we need to spend instructions and time somewhere else (cell lists, packing unpacking etc)

# CTC-ICHIP



# COMPARISON WITH STATE-OF-THE-ART

---

Reference work:

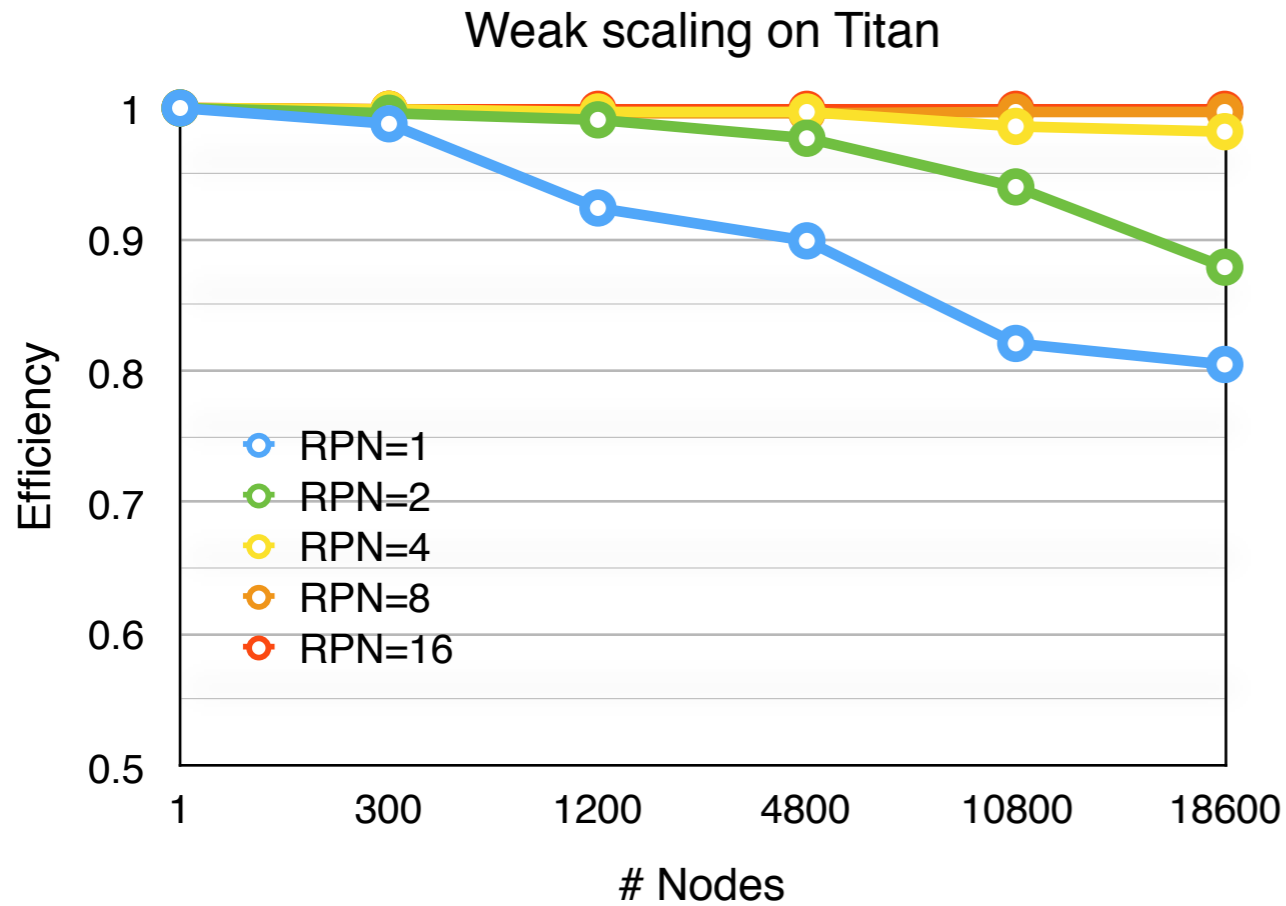
- $7.3 * 10^{10}$  unknown/s on TSUBAME 2
  - $1.0 * 10^{12}$  unknown/s on TITAN assuming perfect scaling
- **uDeviceX: 7.3X** higher throughput



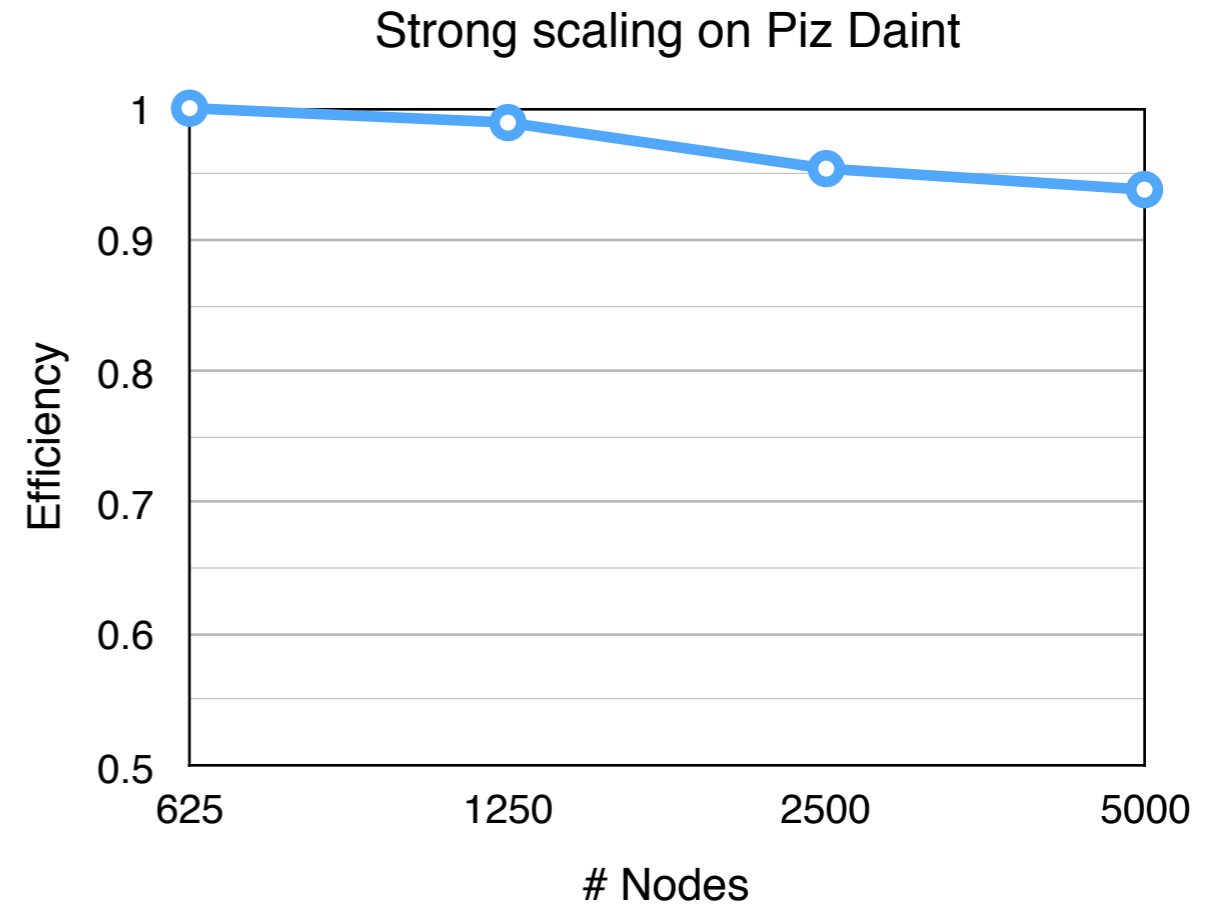
# SCALABILITY

Do we scale? YES!  
Some teams compete for SCALABILITY.

> **98% weak scaling** on the 18k nodes of Titan



**94% strong scaling**  
from 625 to 5000  
nodes of Piz Daint

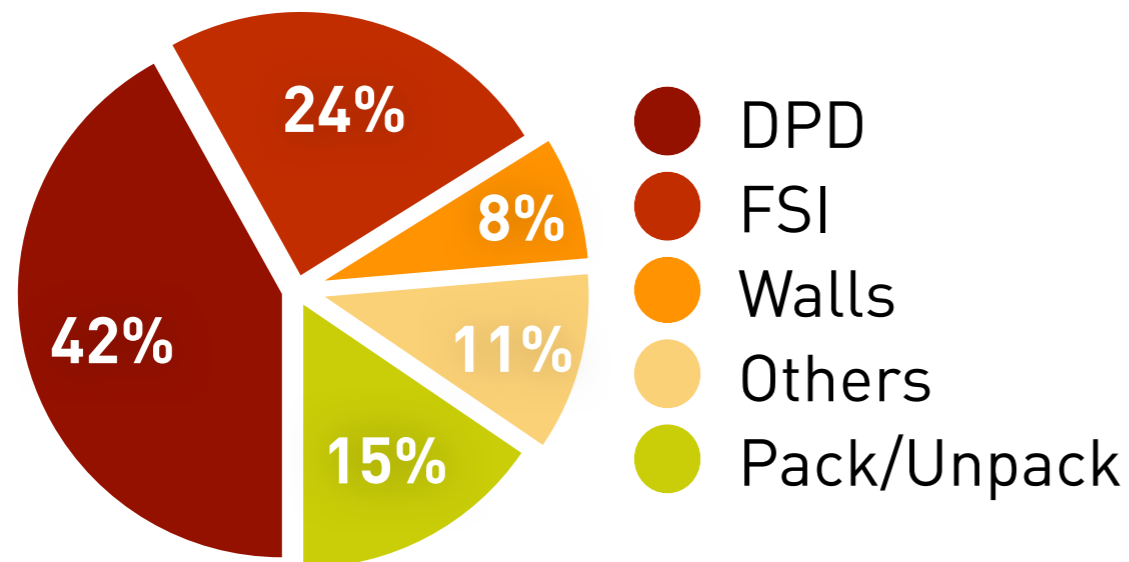




# FURTHER ANALYSIS

---

- **CUDA kernels:** 85% of the GPU time



Kernel	IPC	% GPU peak
<b>DPD</b>	2.7	45%
<b>FSI</b>	2.5	43%
<b>Walls</b>	3.2	52%

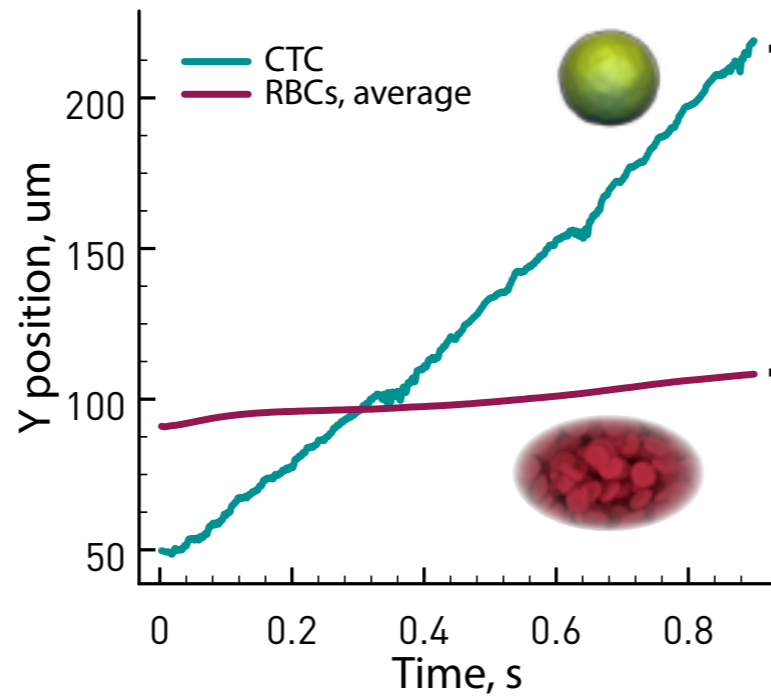
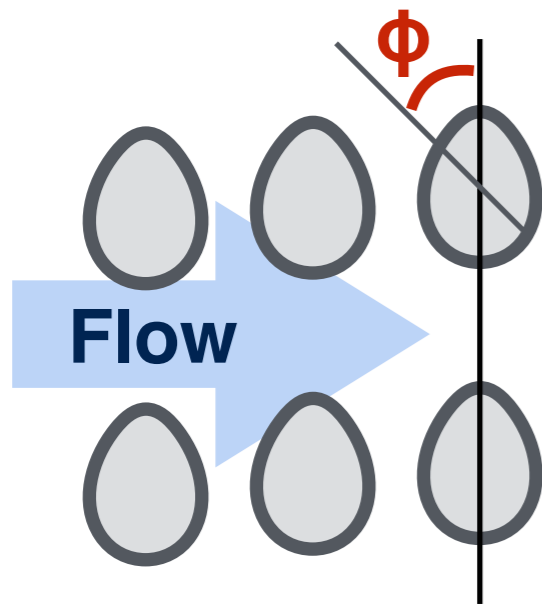
- **Time scales separation**

- 5.5x gain in time-to-solution

- Increased GPU utilization: IPC = 40% of the nominal peak performance

# OPTIMIZING CTC-ICHIP

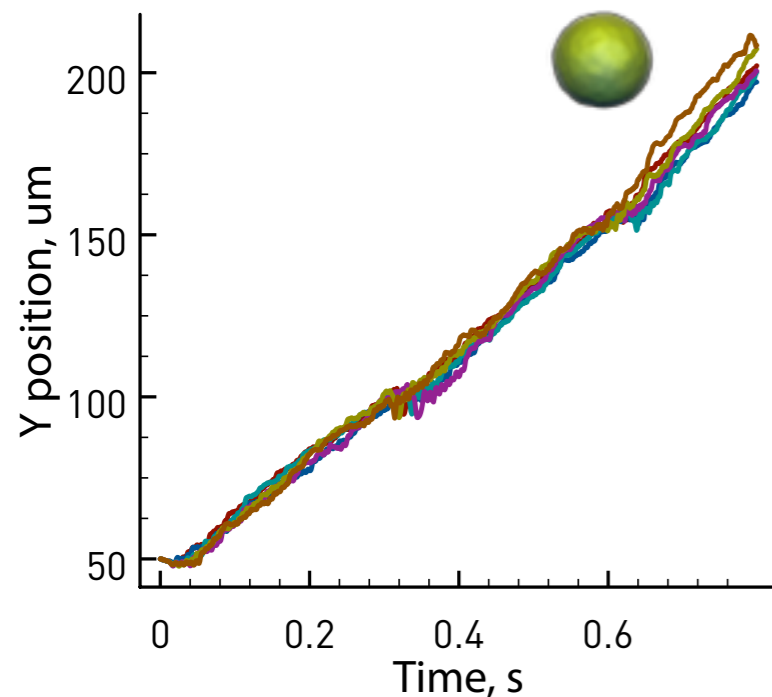
Original CTC-iChip:



Optimizing for separation of CTC and RBCs

Optimizing **angle  $\phi$**

CTC trajectories are similar for different angles



RBCs behave very differently

