

BDMC 2012

Workshop on Big Data
Management in Clouds



Turning Big data into knowledge

Techniques and Tools for Parallel Computing on Online Data Streams in Systems Biology and Epidemiology

Marco Aldinucci
University of Torino, Italy



BioBITs



UNIVERSITÀ
DEGLI STUDI
DI TORINO
ALMA UNIVERSITAS
TAURINENSIS



Outline

- Stochastic Formal System Biology
- From Distributed to Multicore and back
- On programming models
 - FastFlow
- The CWC parallel simulator for sys bio
- Some preliminary results

Formal synthetic system biology



System Biology & Gillespie's algorithm

- Traditionally studied with continuous ordinary differential equations (ODE)
 - bulk reactions, i.e. average behaviour
- Gillespie algorithm: discrete and stochastic simulation of a system via explicit simulation of each reaction
- Gillespie realization represents a random walk that exactly represents the distribution of the master equation (i.e. ODEs)
 - under some hypothesis

Gillespie's algorithm [77]

1. **Initialization:** Initialize the number of molecules in the system, reactions constants, and random number generators.
2. **Monte Carlo step:** Generate random numbers to determine the next reaction to occur as well as the time interval. The probability of a given reaction to be chosen is proportional to the number of substrate molecules.
3. **Update:** Increase the time step by the randomly generated time in Step 2. Update the molecule count based on the reaction that occurred.
4. **Iterate:** goto Step 2 unless the number of reactants is zero or the simulation time has been exceeded.

Increasingly popular approach

- Sometime more informative than (ODE)
 - multi-stability, divergent or rare behaviours, peaks, ...
 - multi-scale systems
 - e.g. deriving macro behaviour from micro
- More computational demanding
 - much more, especially in motivating cases

Increasingly popular approach

- Bio-PEPA [Hillston, Ciocchetta]
- SPiM [Cardelli, Phillips]
- Stochastic Pi [Priami]
- Stochkit [Petzold]
- Spatial Pi [Uhrmacher]
- Calculus of Wrapped Components [our own]
 - kinetics: mass-action, Michaelis–Menten, Hill ...
- ...

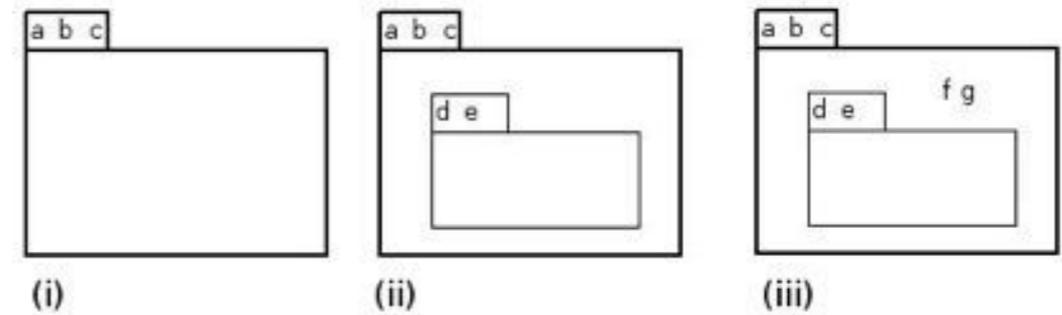
The Calculus of Wrapped Compartments (CWC)

Examples of SCWC terms

A **term** is intended to represent a biological system. A *term* is built by means of the **compartment** constructor, $(- \mid -)$, from a set \mathcal{E} of *atomic elements*, ranged over by a, b, c, d . A **simple term** is defined as:

$$t ::= a \mid (\bar{a} \mid \bar{t})$$

We write \bar{t} to denote a (possibly empty) multiset of simple terms $t_1 \dots t_n$. Similarly, with \bar{a} we denote a (possibly empty) multiset of atoms.



- (i) represents $(a \ b \ c \mid \bullet)$;
- (ii) represents $(a \ b \ c \mid (d \ e \mid \bullet))$;
- (iii) represents $(a \ b \ c \mid (d \ e \mid \bullet) \ f \ g)$.

Dynamics of SCWC

Stochastic Rules

Rewrite rules are defined as pairs of terms, in which the left term characterizes the portion of the system in which the event modelled by the rule can occur, and the right one describes how that portion of the system is changed by the event.

Biomolecular Event	Examples of CWC Rewrite Rules
State change	$a \mapsto b$
Complexation	$a \ b \mapsto c$
Catalyzed membrane crossing	$a \ (b \ x \mid y) \mapsto (b \ x \mid a \ y)$ $(b \ x \mid a \ y) \mapsto a \ (b \ x \mid y)$

Rules are decorated with a **rate** (speed of the reaction).

A **Stochastic Rewrite Rule**, R , is denoted by $P \xrightarrow{k} P'$.

The stochastic semantics is given by transitions between terms labeled with the rule applied, R , and a transition rate depending on the rate of rule R :

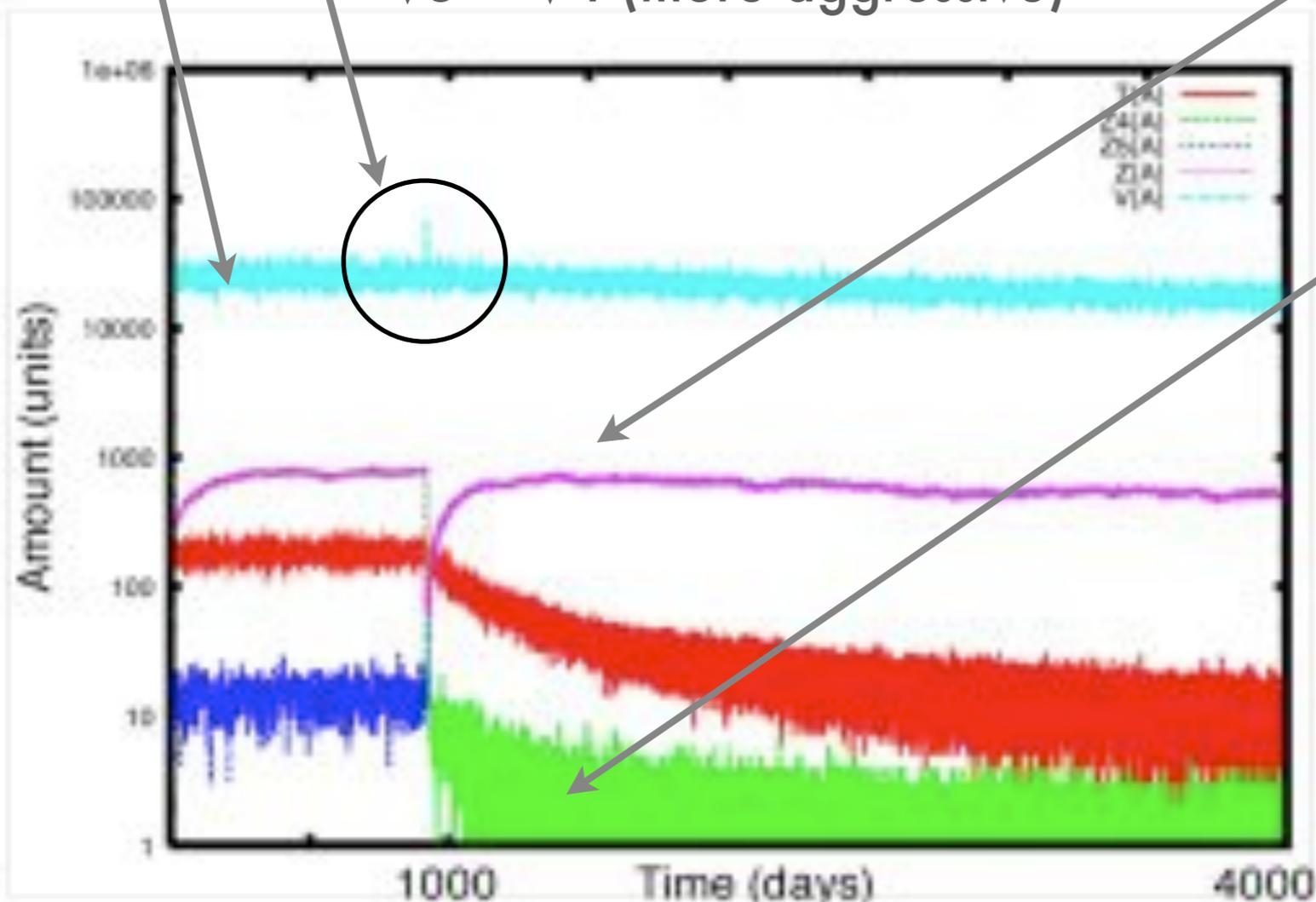
$$\bar{t} \xrightarrow{R, k \times p} \bar{t}'$$

where R is $P \xrightarrow{k} P'$, and p is the number of different ways in which the pattern P may match \bar{t} ($\bar{t} = C[P\sigma]$) and such that $\bar{t}' = C[P'\sigma]$ for some context C and variable instantiation σ .

Ex: HIV and immune response (progression to AIDS)

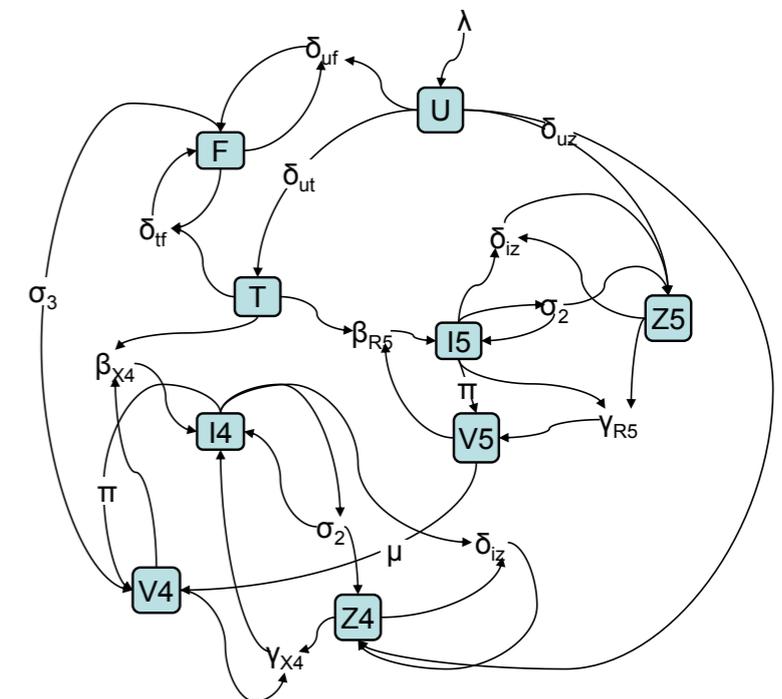
V virus, all phenotypes (V4,V5)

the spike suggests the mutation
V5 → V4 (more aggressive)



immune response $Z=Z4+Z5$
remain stable (but for the peak).

now $Z4$ decrease and
 $Z5$ increase
i.e. HIV is turning into
AIDS more rapidly



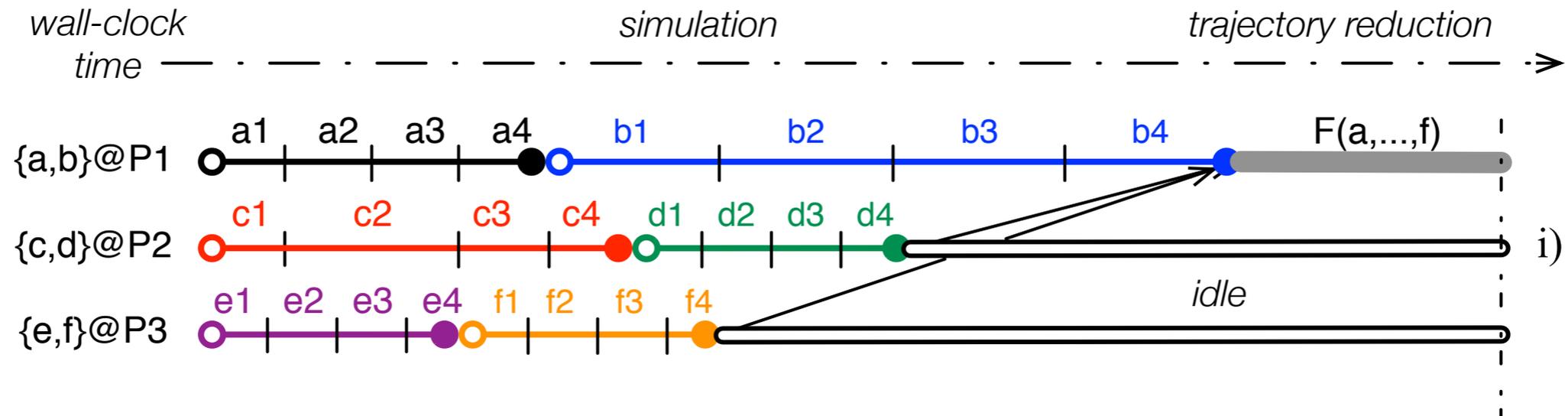
- Peaks are informative events,
 - virus mutation triggers AIDS progression
 - hardly detected with ODEs
- high resolution required to detect spikes,
 - each trajectory can be over 6G Bytes of data
- and thousands of trajectories are needed
 - compute everything, save everything, move and join all data, analyse all data, then get first results
 - often to discover parameters are wrong ...



- It is Monte Carlo,
 - well understood
 - easy to parallelise on different trajectories
- it is Monte Carlo w Markov Chains models (CTMC)
 - single trajectory: no parallelisation without relaxation
 - compute time \neq simulation time
 - compute time for different trajectories heavily unbalanced
 - fast reactions and slow reactions, some not interesting (e.g. water-steam-water)



Unbalancing + filtering

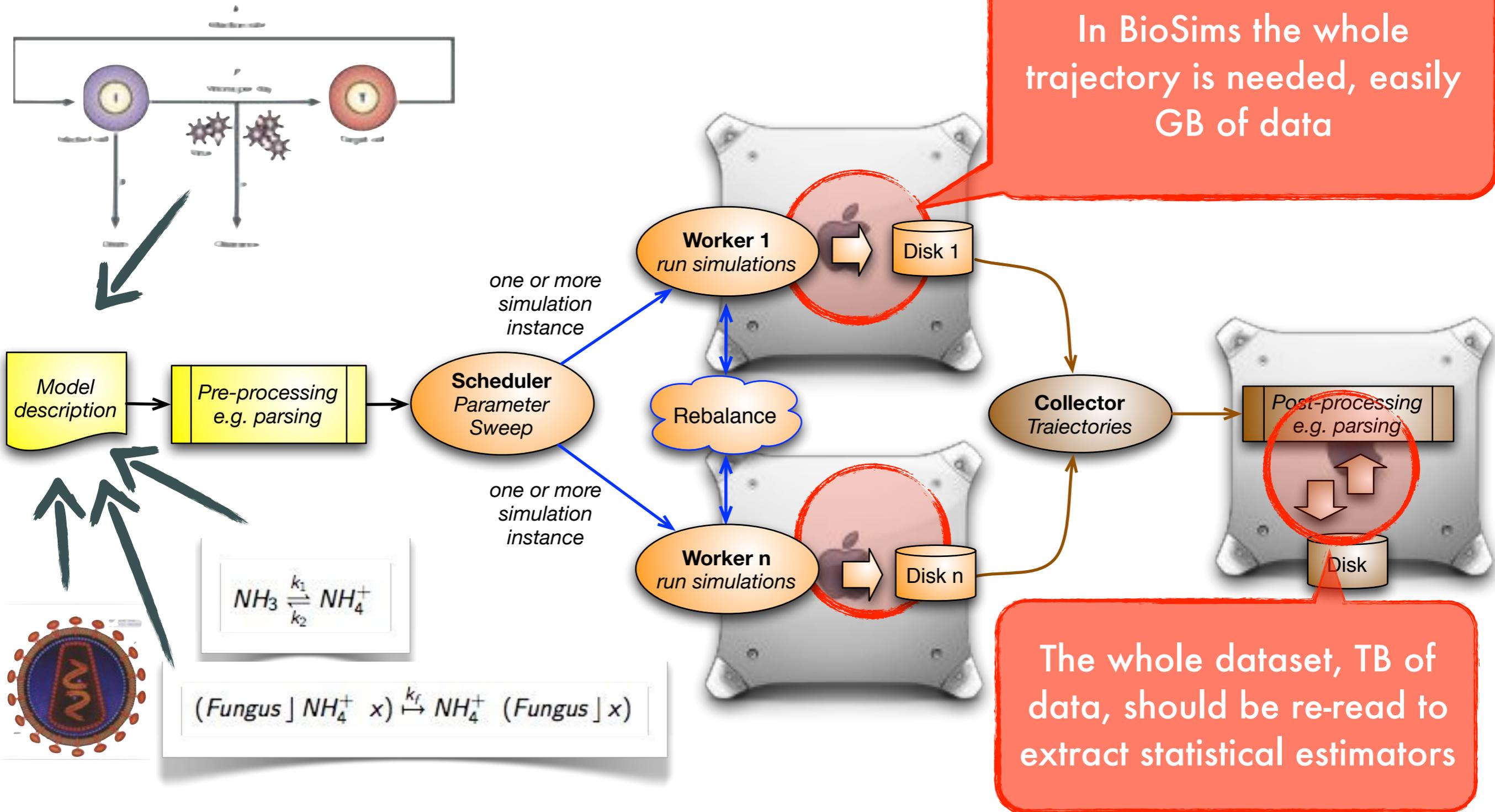


- few trajectories (e.g. the interesting ones) can significantly delay the completion of others
 - over-provisioning don't help that much
 - simulated time moves at different pace w.r.t. wall-clock time
 - data joining from different trajectories should be aligned at the same simulation time

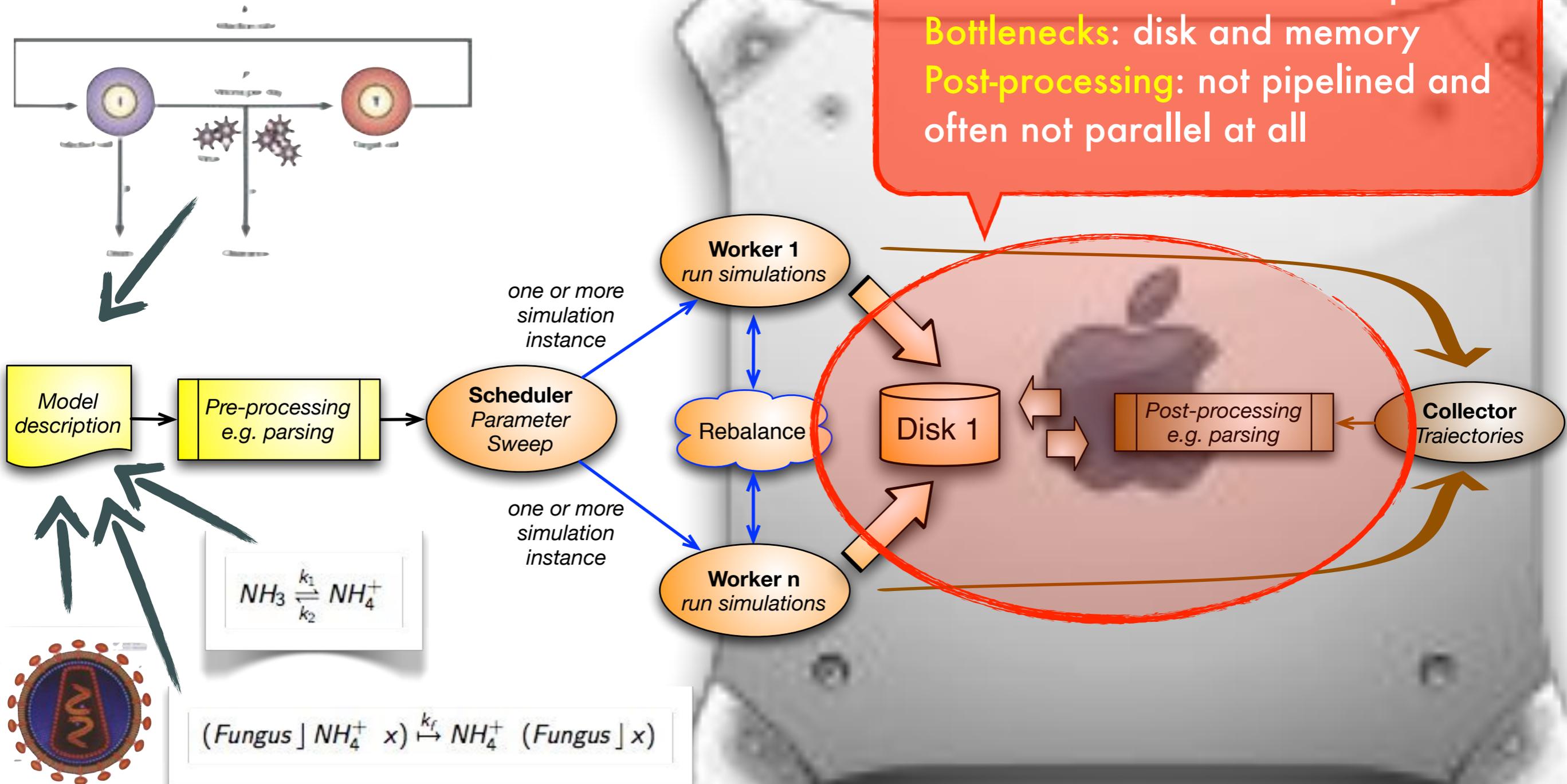
- It is Monte Carlo,
 - well understood
 - easy to parallelise on different trajectories
- it is Monte Carlo w Markov Chains models (CTMC)
 - single trajectory: no parallelisation without relaxation
 - compute time \neq simulation time
 - compute time for different trajectories heavily unbalanced
 - fast reactions and slow reactions, some not interesting (e.g. water-steam-water)
- it is Monte Carlo AND data analysis
 - data is big, analysis can be very expensive and **it typically starts after the simulation**
 - the whole workflow is perceived too “slow” by bio-scientists to be really useful



From Distributed to Multicore and back



From Distributed to Multicore and back



From Distributed to Multicore and back

- Multi Carlo sims for Bio are I/O-bound
 - Sampling reduce I/O traffic but worsen precision and analysis of “strange” dynamics (spikes, diversion from average, etc.), which observation motivates stochastic analysis (ODEs)
- Data analysis is also I/O-bound
 - if approached is a “post-processing” fashion, data should be retrieved from the disks
- The porting of distributed solution “as is” on multicore is going insist on weakness of multicore architectures
 - Memory wall, I/O, disk
 - SIMD / GPGPUs do not change the analysis substantially

From Distributed to Multicore and back

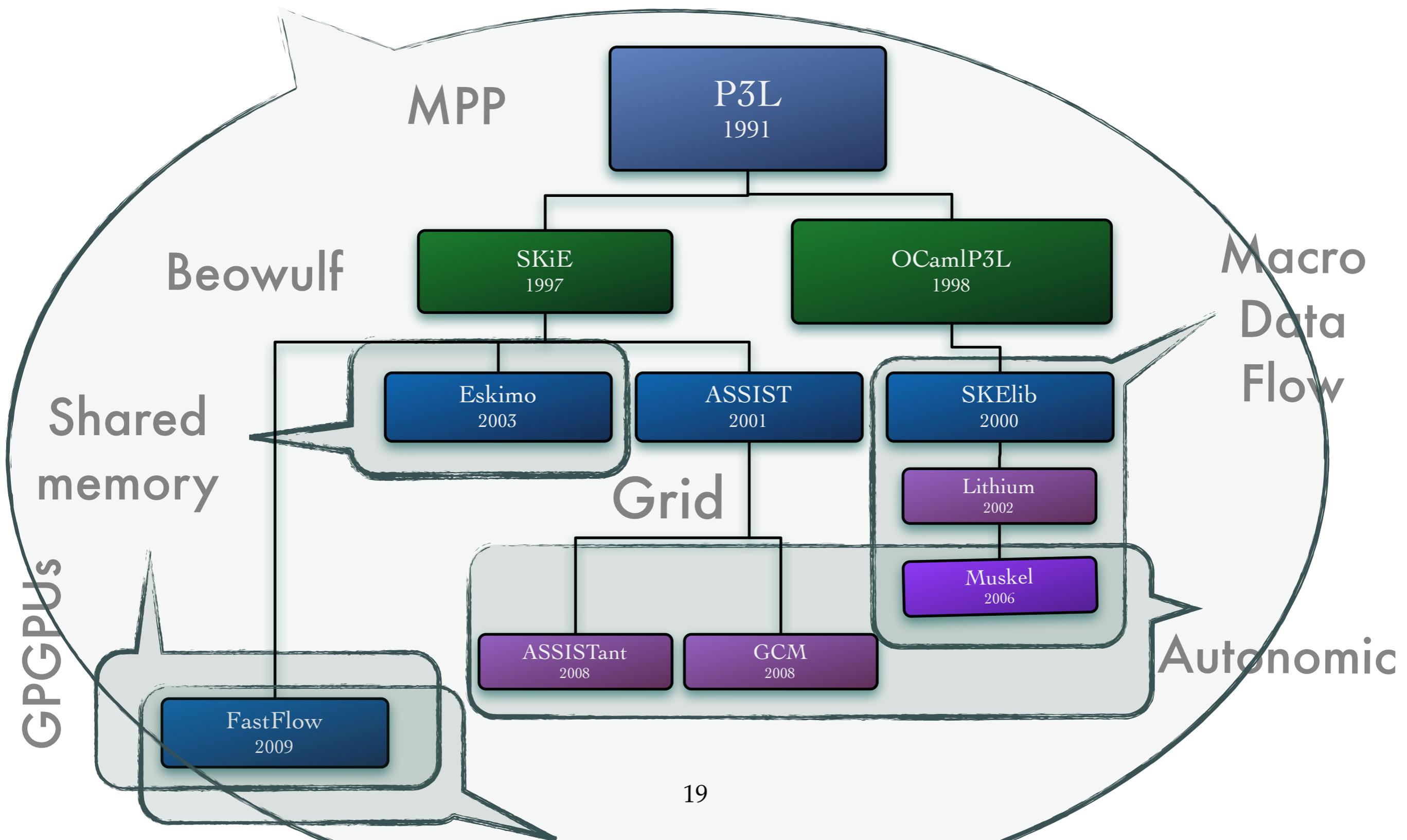
- The same arguments holds on distributed, grids, and clouds as soon as the workflow is considered as a whole
 - simulation, data collection and merging, analysis
- Rationale
 - **Manage data as stream, compute online**
 - May require more computing and less bandwidth
 - Computation should be designed to be pipelined
 - **Establish fast data paths across cores/hosts**
 - **Avoid low-level concurrency management**
 - Portability, performance, portability of performance, maintenance, porting from sequential

Quotes from P. Beckman EuroPar 11

“exascale” keynote

- *coarse grain concurrency is nearly exhausted*
- *it is not about Flops, it is about data movement*
- *programming systems should be designed to support fast data movement and enforce locality*
- *shared-memory & inter-socket messaging*
- *we need a programming model*
 - *a computer language is not a computing model*
 - *a library is not a computing model*
- *we need a efficient and compositional run-time*

Patterns/skeletons & streams



Clouds, clusters of multi-cores and many-cores

it is not about Flops, it is about data movement, we need compositional run-time

- Streams
 - focus on data movements at the prog model level
 - clear semantics
 - support compositionally and also locality
 - the latter is a bit more counter-intuitive
- High-level programming
 - e.g. patterns
- Patterns + streams
 - can be implemented efficiently on both multi-core, distributed, and both



FastFlow

<http://mc-fastflow.sourceforge.net/>

FastFlow (multicore)

Applications on multicore, many-core
Efficient and portable - designed with high-level patterns

FastFlow

Streaming network patterns
Skeletons: pipeline, map farm, reduce, D&C, ...

Arbitrary streaming networks
Lock-free SPSC/MPMC queues + FF nodes

Simple streaming networks
Lock-free SPSC queues + threading model

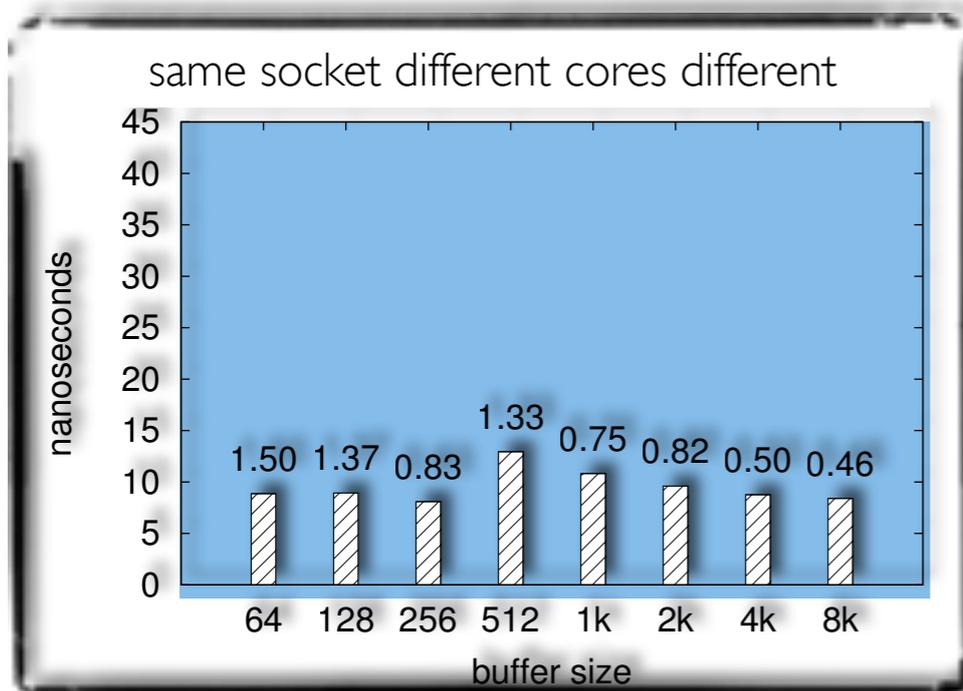
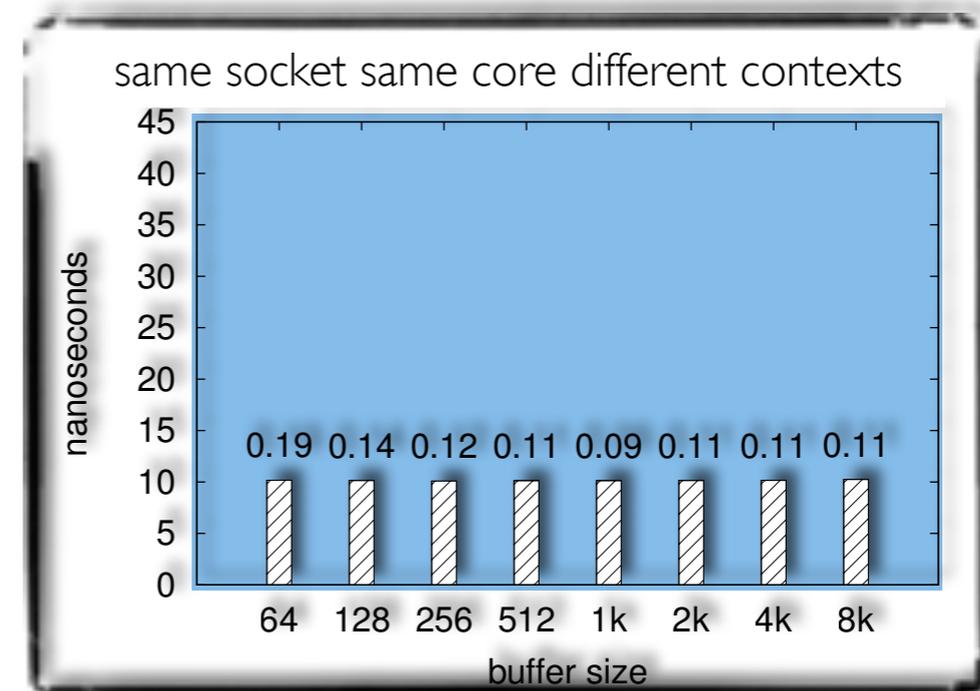
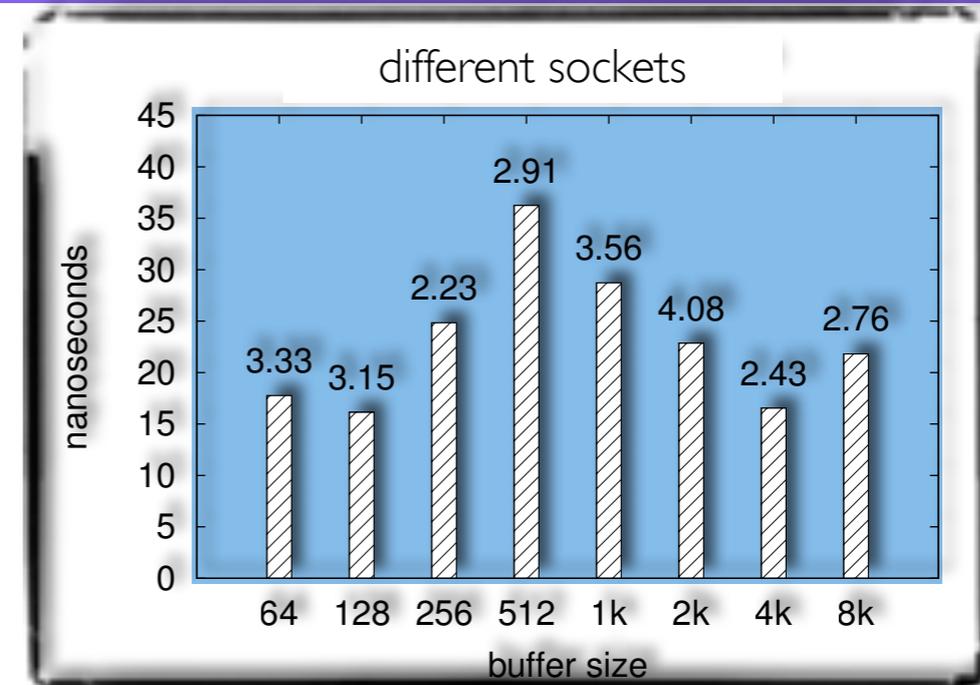
Multicore and manycore
SMP: cc-UMA & cc-NUMA

Layer 1: Simple streaming networks

4 sockets x 8 core x 2 contexts

Xeon E7-4820 @2.0GHz Sandy Bridge
18MB L3 shared cache, 256K L2

MPI is ~190 ns at best
(D.K. Panda)



Layer 1: Simple streaming networks

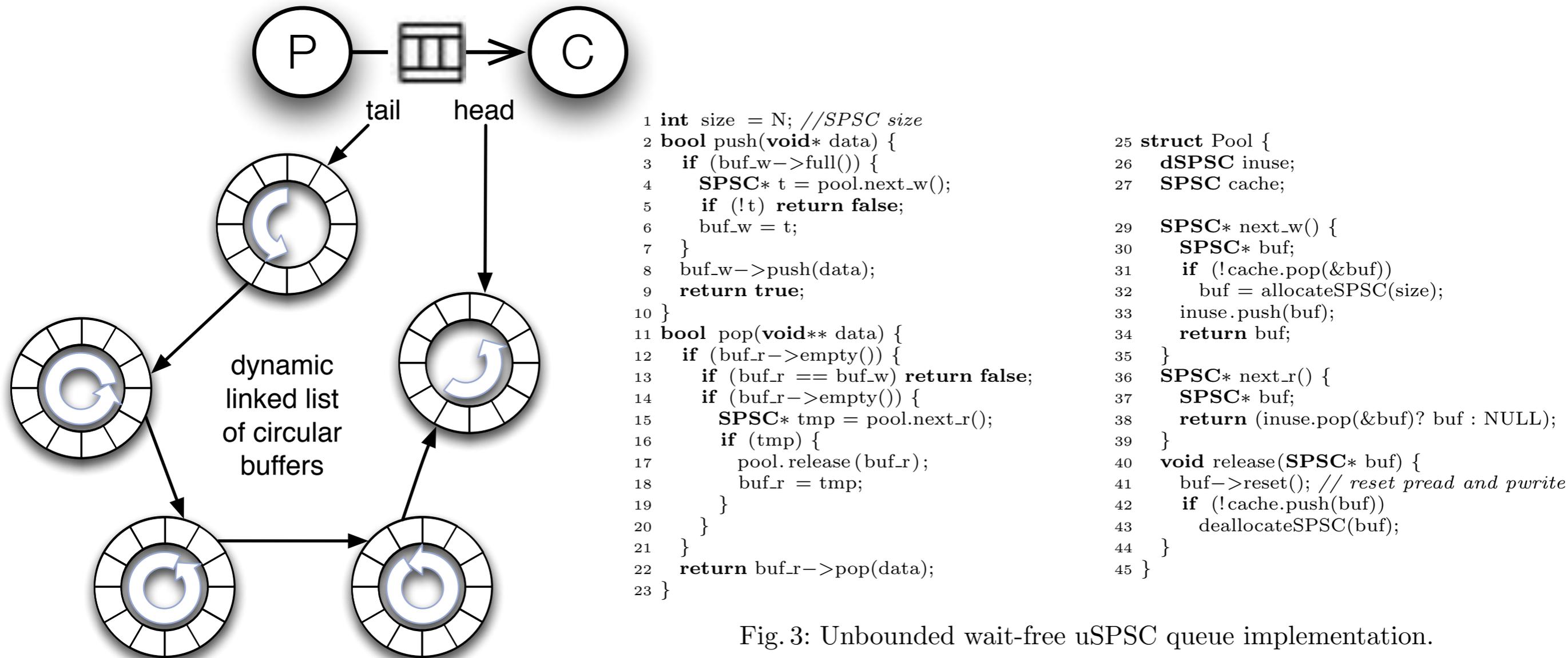


Fig. 3: Unbounded wait-free uSPSC queue implementation.

M. Aldinucci, S. Campa, M. Danelutto, M. Torquati. An Efficient Synchronisation Mechanism for Multi-Core Systems.

EuroPar 2012.

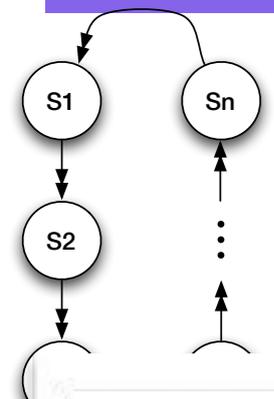
Wed 29 Aug - B3 multicore 14.30-16.00

Layer 1: Simple streaming networks

<http://www.1024cores.net/home/technologies/fastflow>

4 sockets x 8 core x 2 contexts

Xeon E7-4820 @2.0GHz Sandy Bridge
18MB L3 shared cache, 256K L2



<35 ns irrespectively of the mapping

Speedup

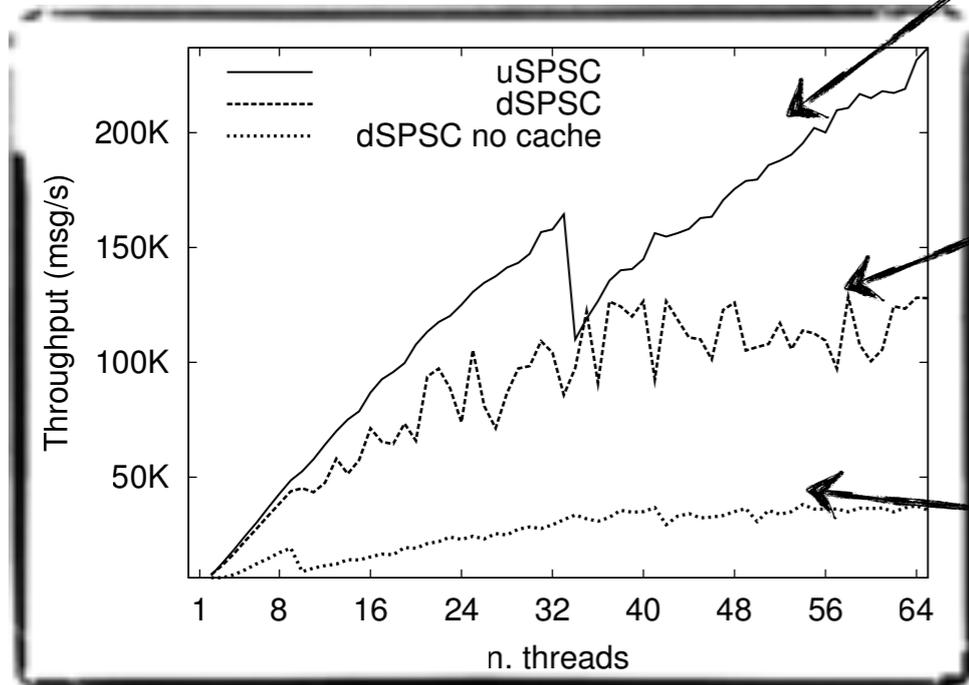
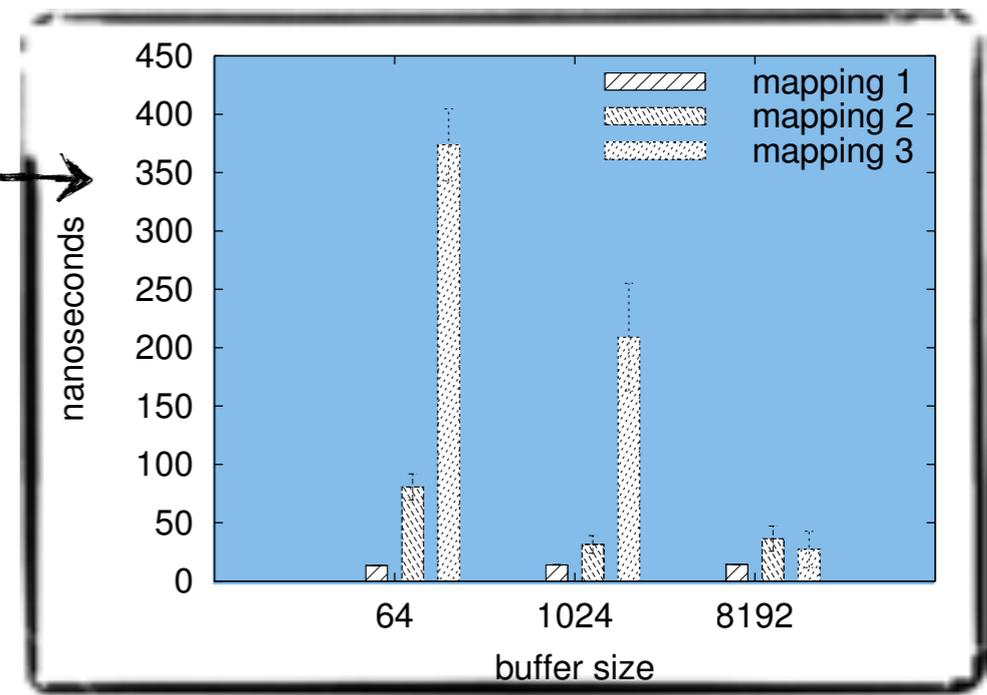
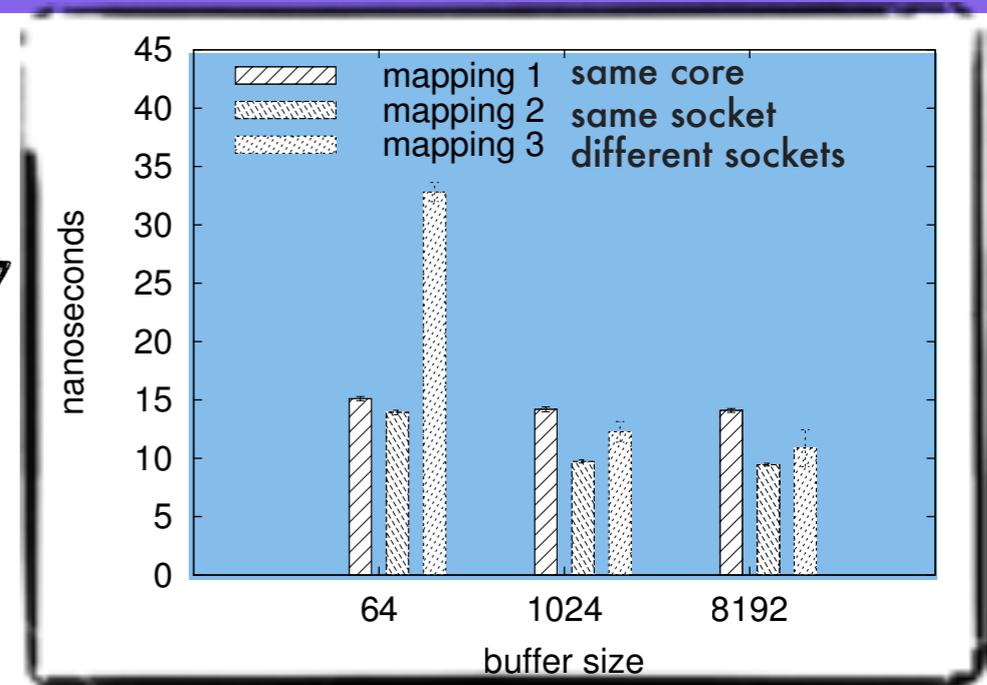
Linked list + circular buffer
FastFlow queue
(our result)

12x faster

Linked list with pooling
Opt Michael-Scott queue
(our result)

20x faster

Linked list w/ dyn alloc
Michael-Scott queue
well-known ~ 400 citations



Layer 3: streaming networks patterns

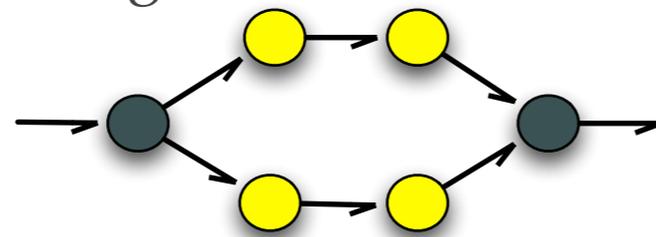
- Composition via C++ template meta-programming

- CPU: Graph composition

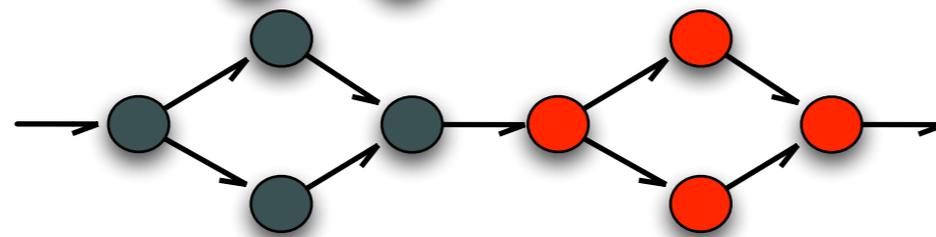
- GPU: CUDA streams

- CPU+GPU: offloading

- `farm{ pipe }`

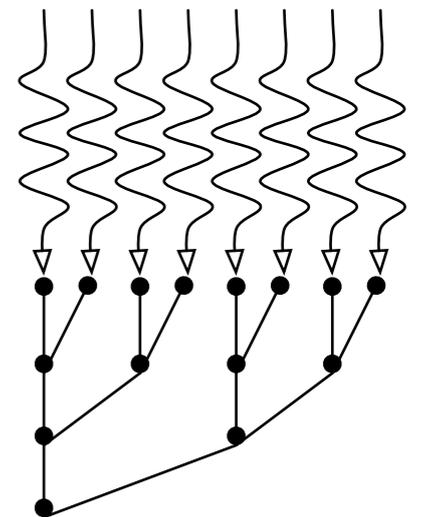
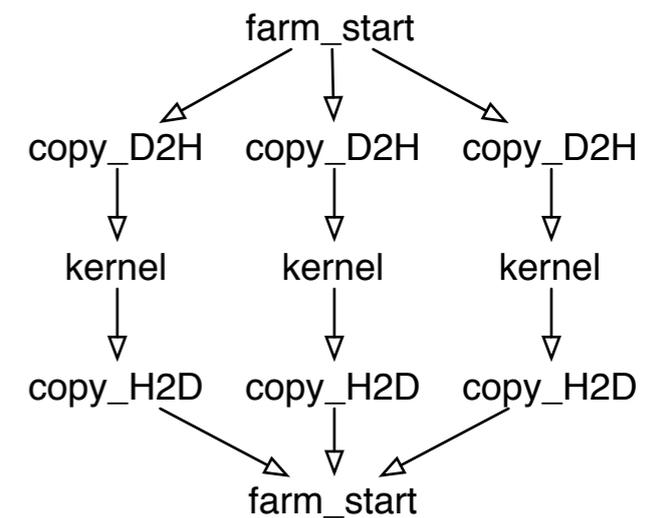


- `pipe(farm, farm)`



- `pipe(map, reduce)`

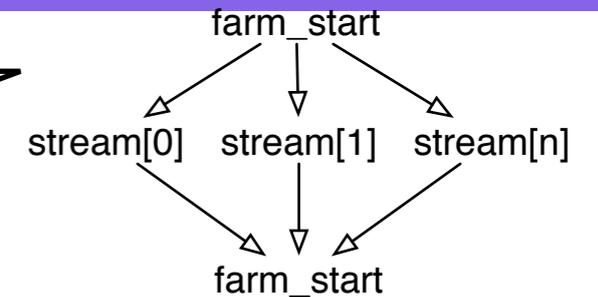
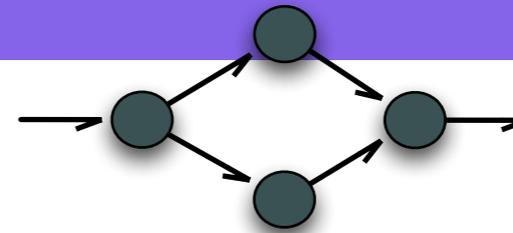
-



Layer 3: streaming networks patterns

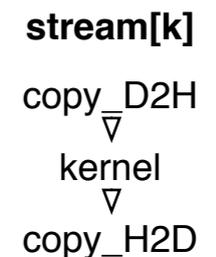
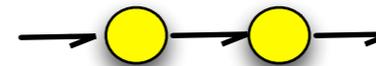
- farm

- on CPU - master-worker - parallelism exploitation
- on GPU - CUDA streams - automatic exploitation of asynch comm



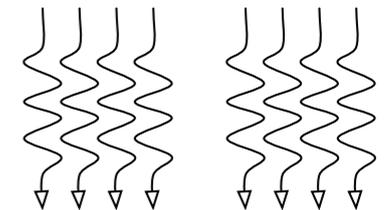
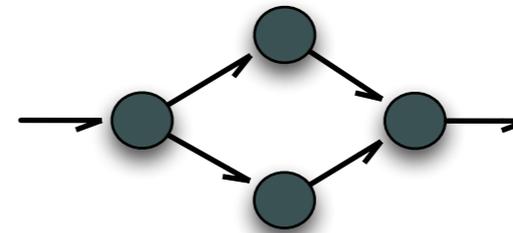
- pipeline

- on CPU - pipeline
- on GPU - sequence of kernel calls or global mem synch



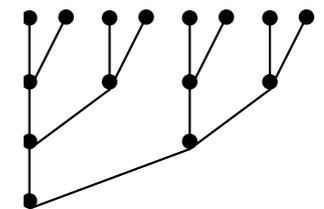
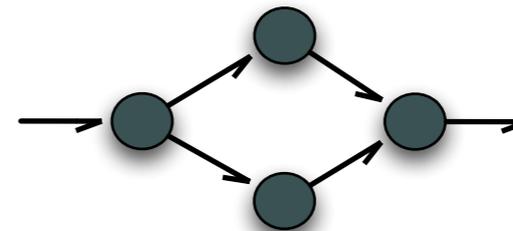
- map

- on CPU - master-worker - parallelism exploitation
- on GPU - CUDA SIMT - parallelism exploitation



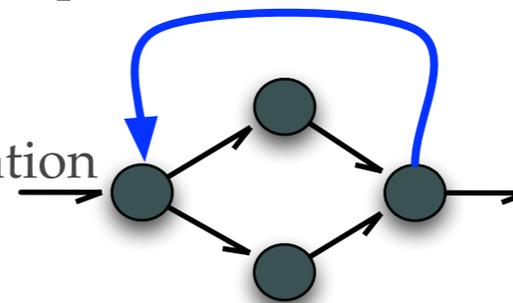
- reduce

- on CPU - master-worker - parallelism exploitation
- on GPU - CUDA SIMT (reduction tree) - parallelism exploitation



- D&C

- on CPU - master-worker with feedback - // exploitation
- on GPU - working on it, maybe loop+farm



Layer 3: streaming networks patterns (easy to port)



+ distributed

Applications on multicore, many core & distributed platforms of multicores
Efficient and portable - designed with high-level patterns

FastFlow

Streaming network patterns

Skeletons: pipeline, map farm, reduce, D&C, ...

Arbitrary streaming networks

Lock-free SPSC/MPMC queues + FF nodes

Arbitrary streaming networks

Collective communications + FF Dnodes

Simple streaming networks

Lock-free SPSC queues + threading model

Simple streaming networks

Zero copy networking + processes model

Multicore and manycore

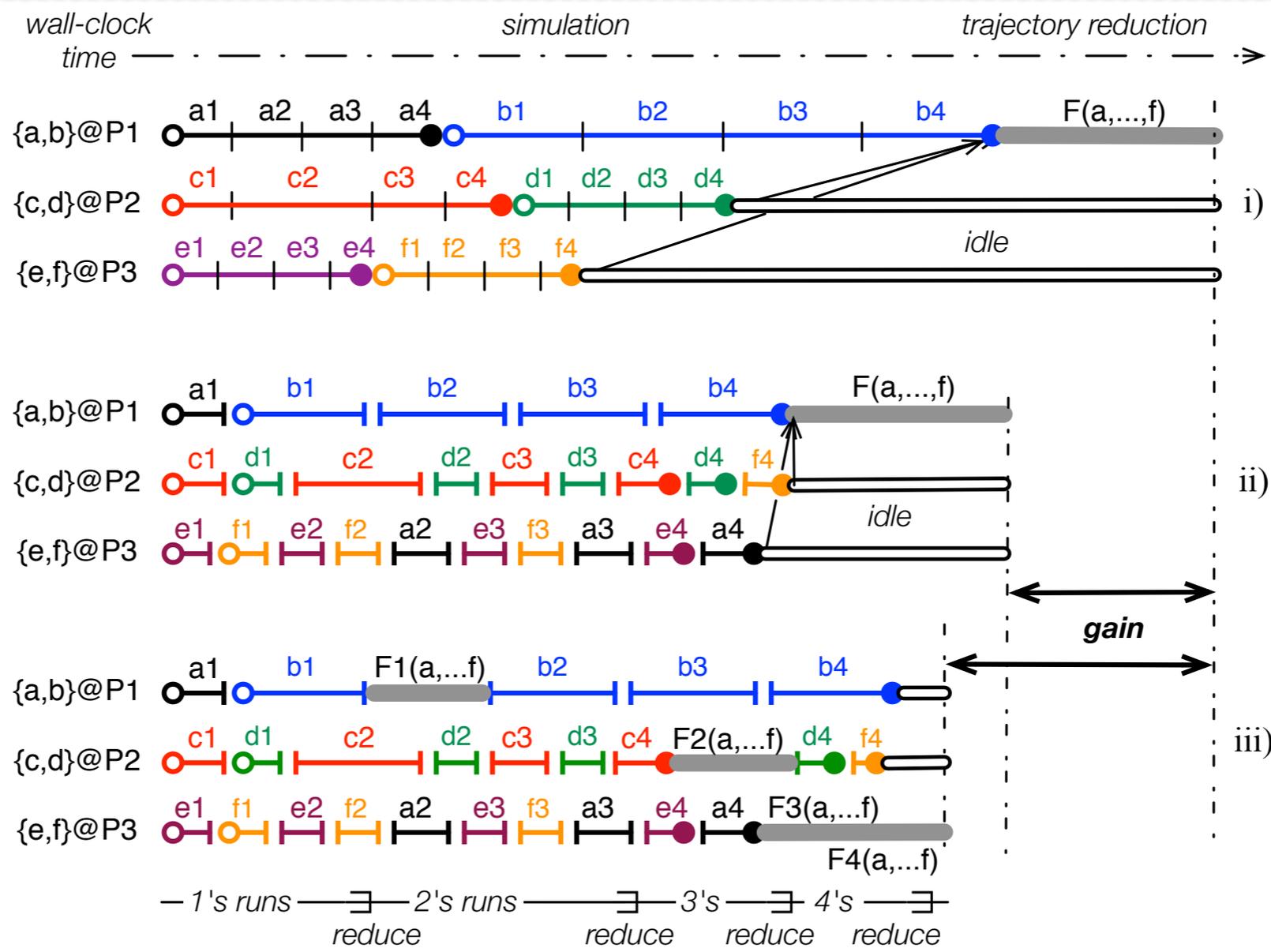
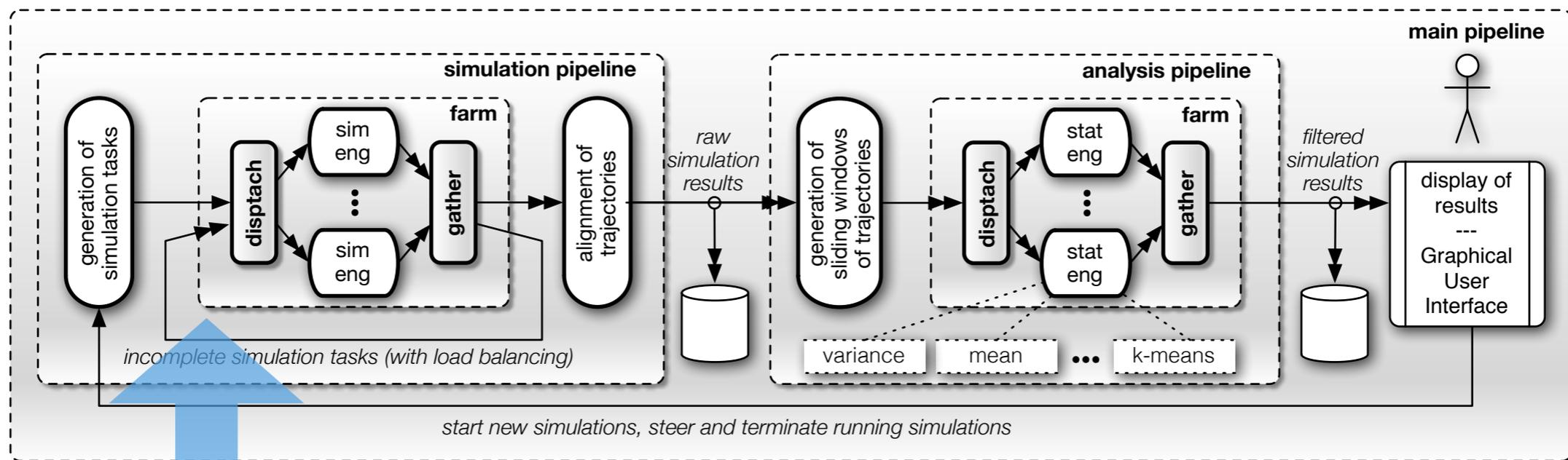
SMP: cc-UMA & cc-NUMA

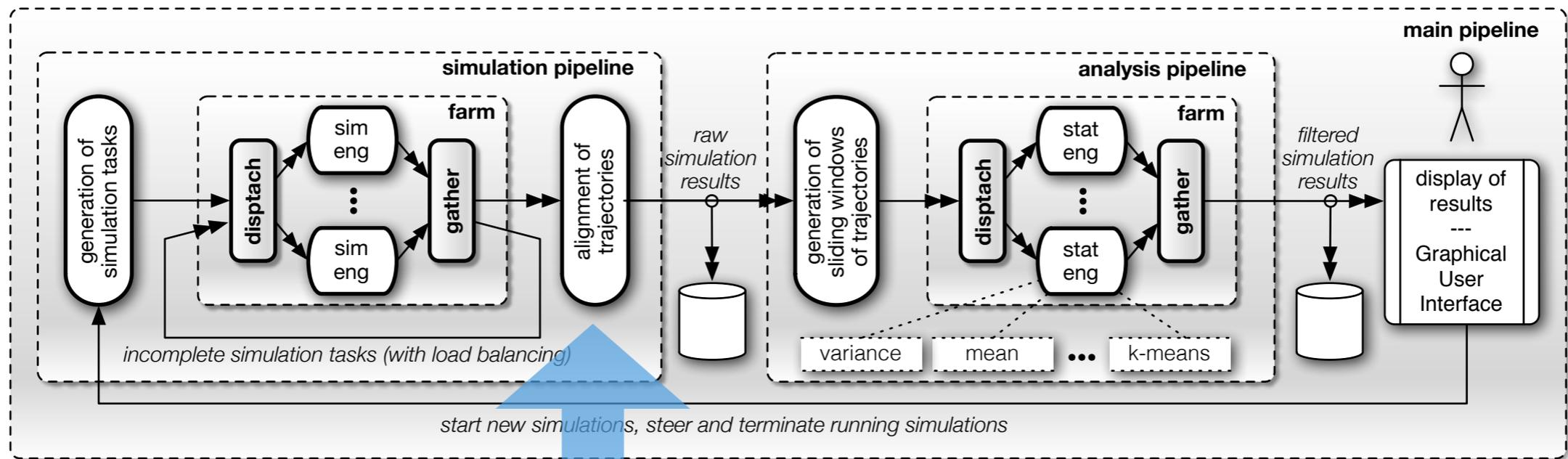
Distributed platforms

Clouds, clusters of SMPs

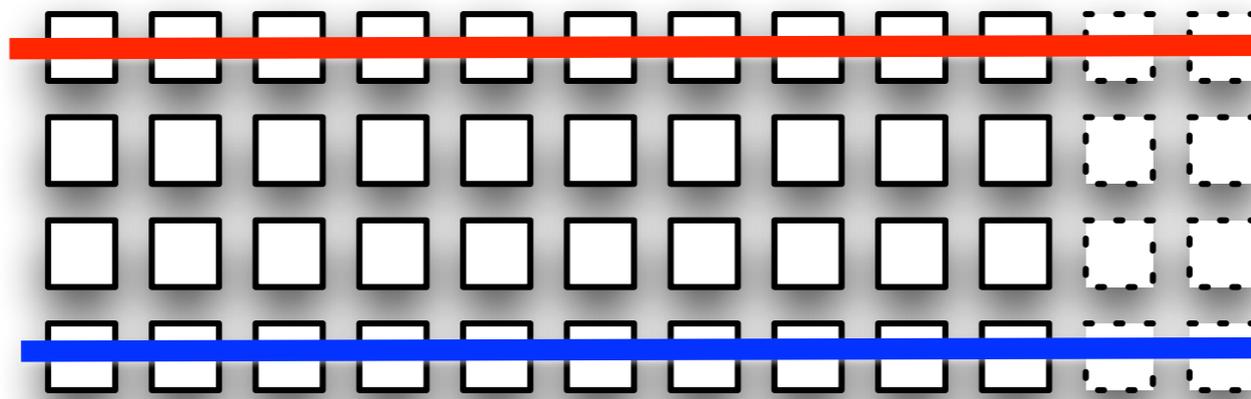
- Generic `ff_node` is subclassed to `ff_dnode`
- `ff_dnode` can support network channels
 - P2P or collective
 - used as frontier node of streaming graph
 - can be used to merge graphs across distributed platforms
- No changes to programming model
 - at least require to “add” stub `ff_dnode`
 - when passing pointers data is serialised
 - serialisation hand-managed (zero-copy, think to Java!)

CWC simulator example





trajectories

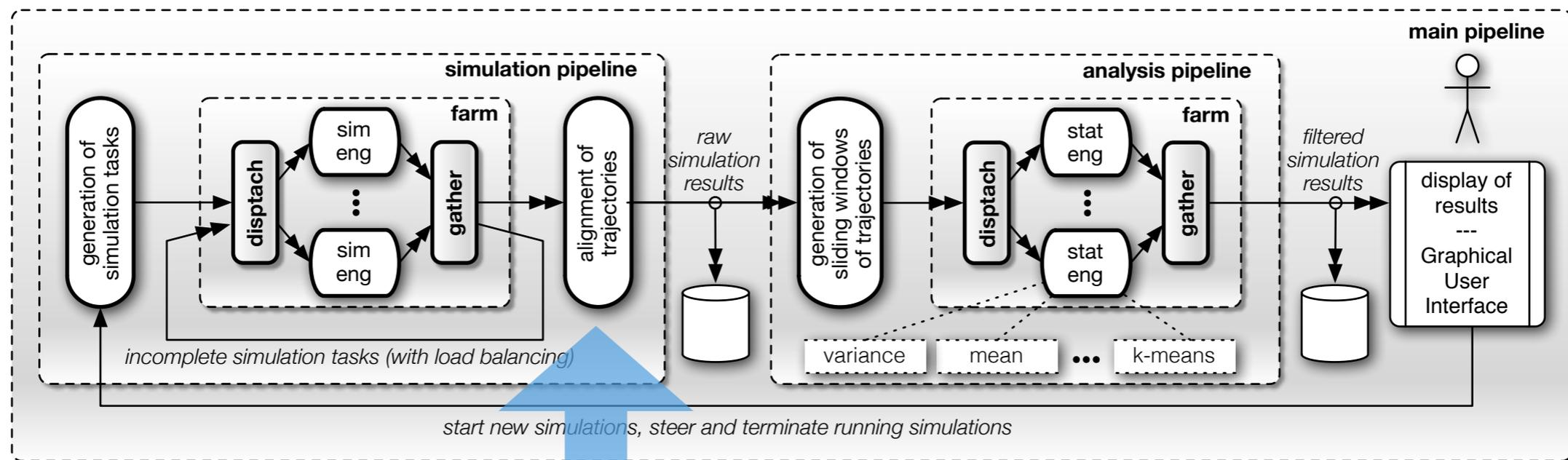


Trajectory 1

⋮

Trajectory n

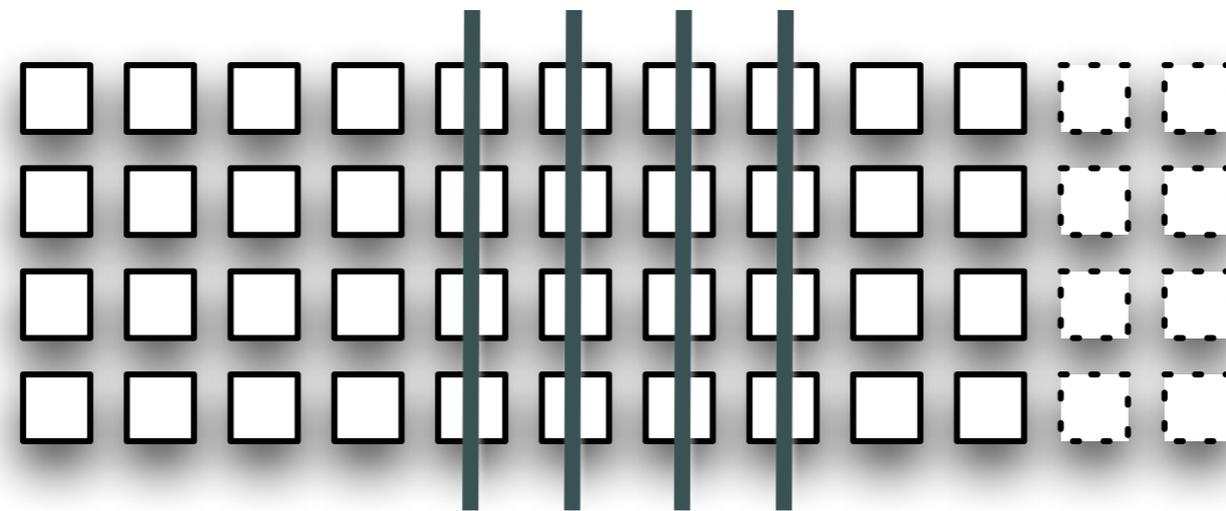




start new simulations, steer and terminate running simulations

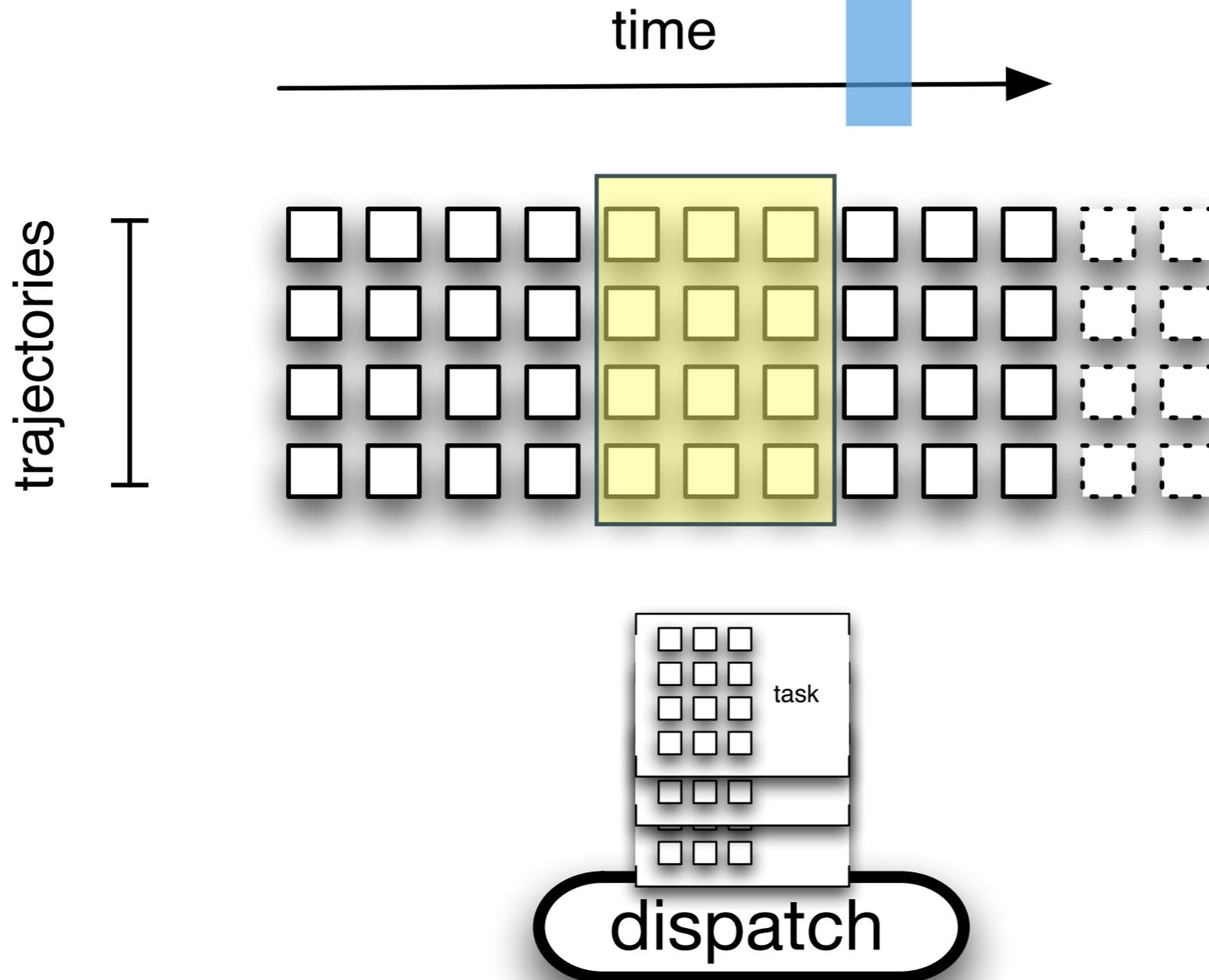
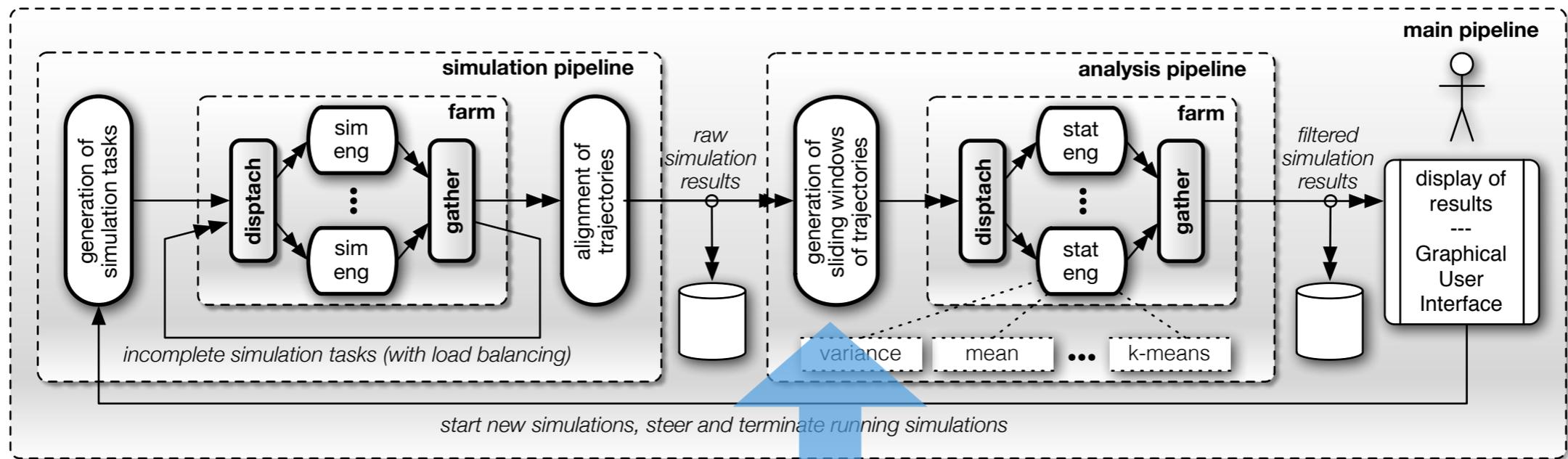
time

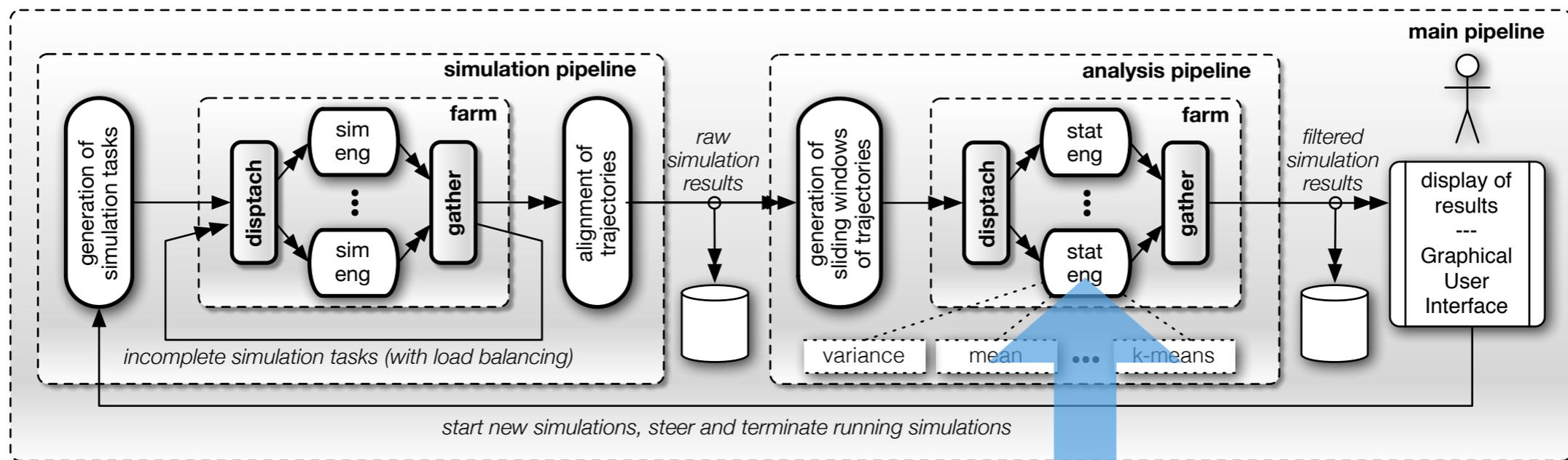
trajectories



simulation-time aligned trajectories





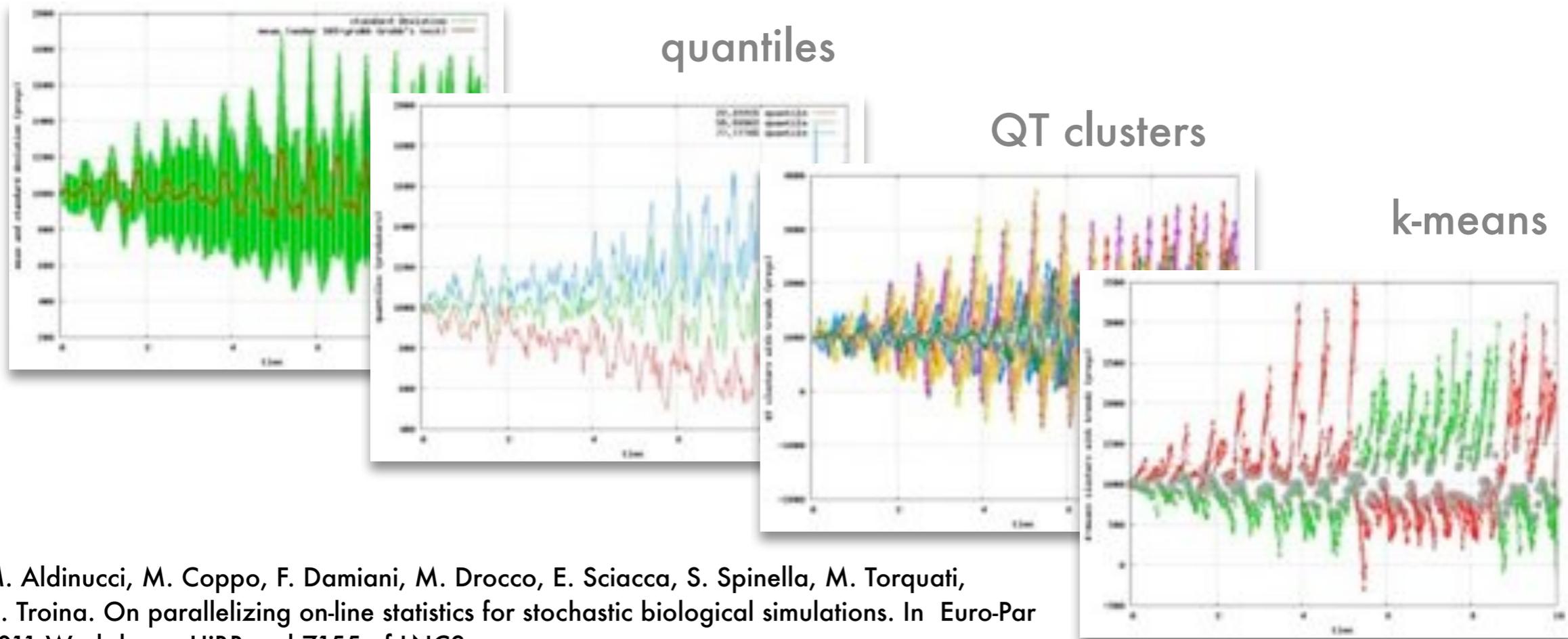


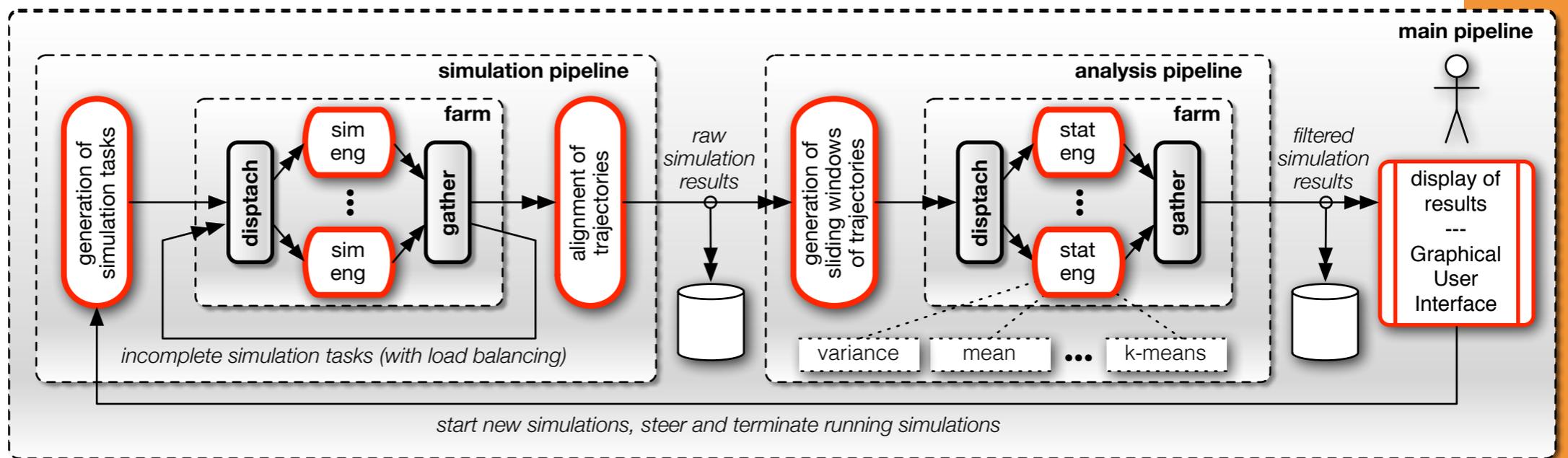
mean variance

quantiles

QT clusters

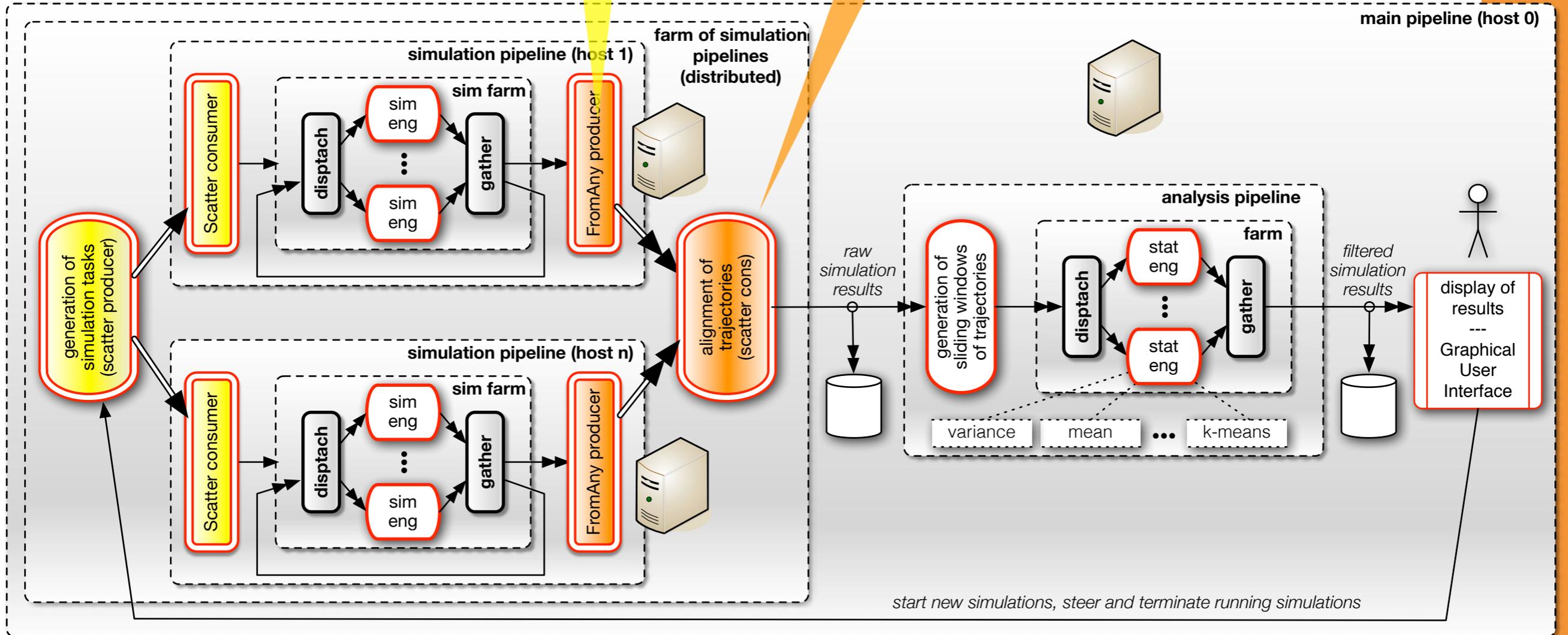
k-means



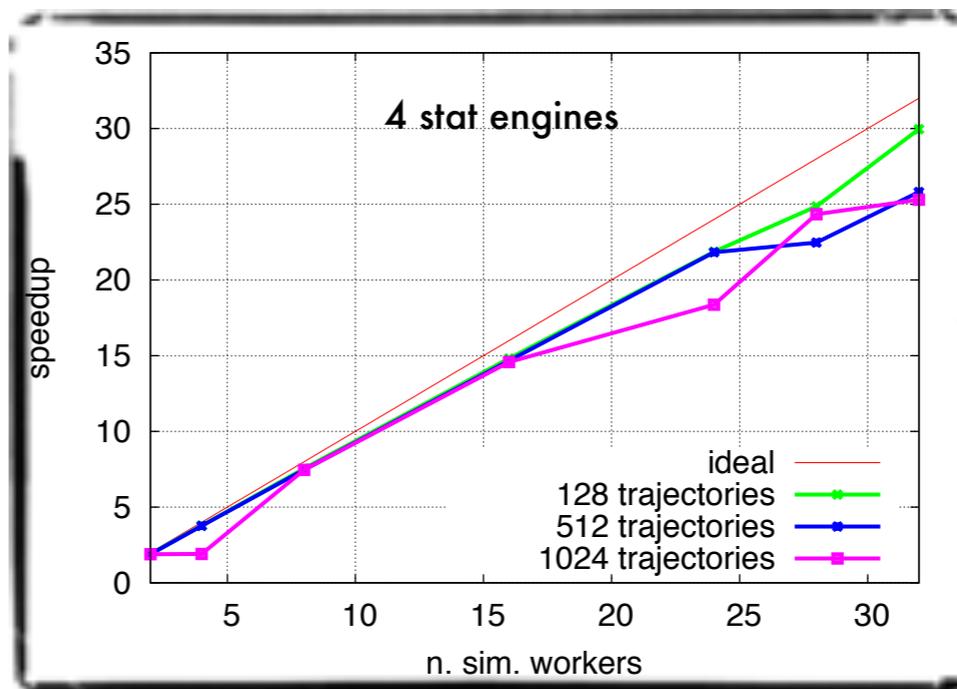


serialisation, coalescing

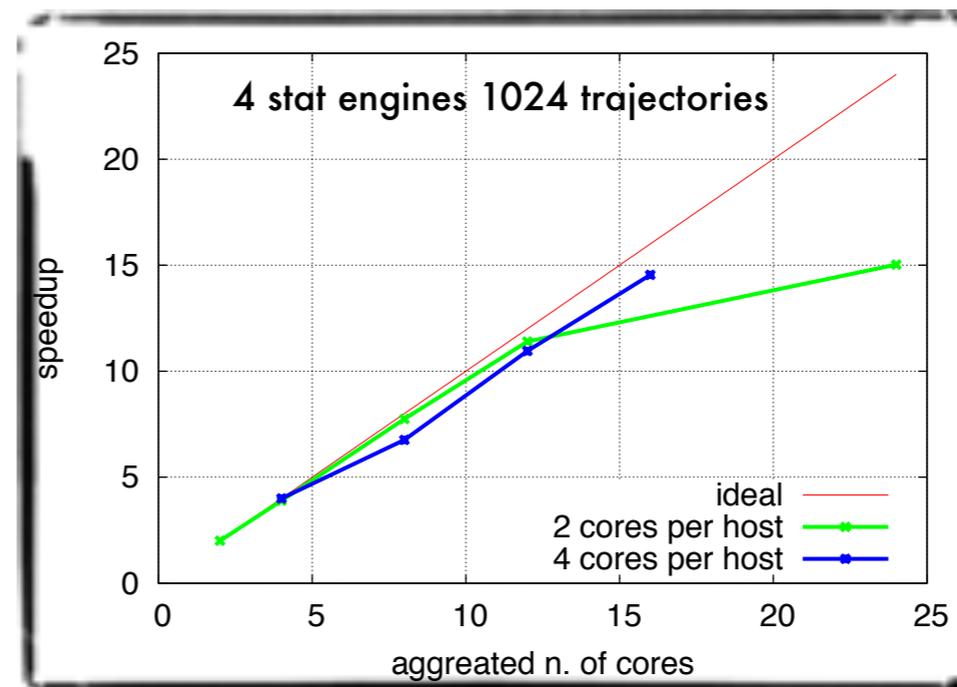
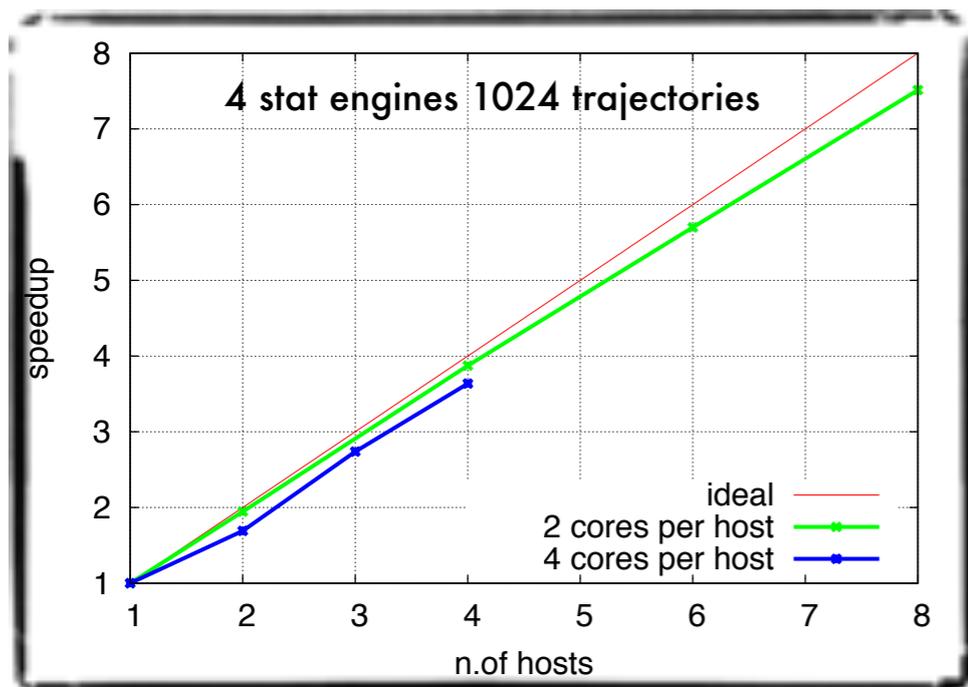
allocation (zero-copy OMQ), de-serialisation



Performance (preliminary)



Intel
Nehalem
32-core



Cluster 8x
Intel Xeon 6-core
Infinband (IPoIB)

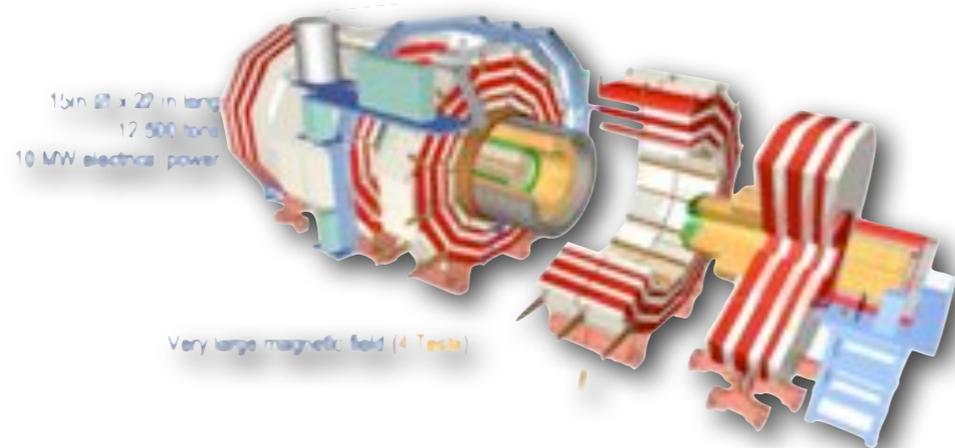
courtesy of
Mellanox

Involved data

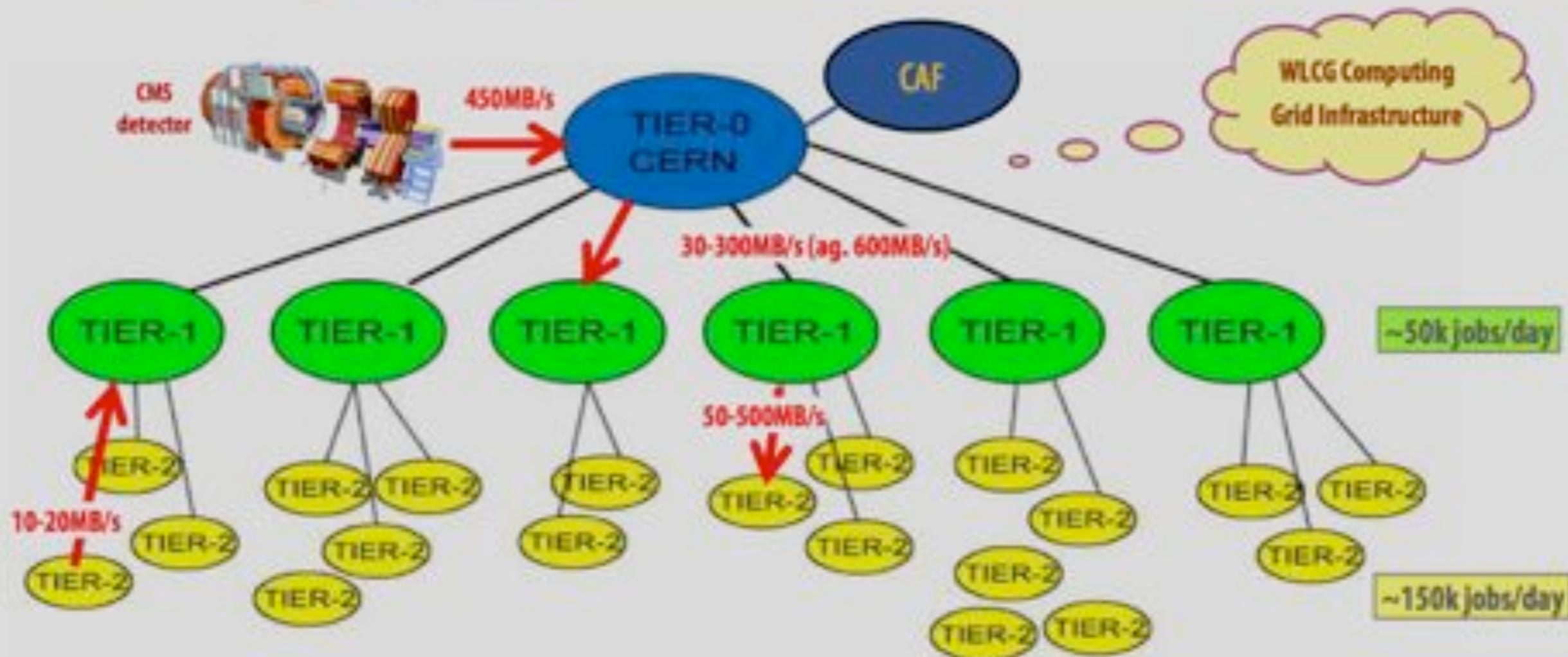
- Simple examples (neurospora, ...)
 - $2-4 \text{ double} * \text{n. of variables} * \text{n. of samples} * \text{n. of trajectories} * \text{cases in sensitivity analysis}$
 - e.g. $4 * 8 * 4 * 1\text{M} * 1\text{k} * 8 \sim 1 \text{ TBytes}$
- HIV 6GB x 1024 trajectories $\sim 6\text{TB}$
- The more observed variables, precision, cases for sensitivity analysis the more data



Innovative Methods for Particle Colliders at the Terascale (2012-2015)



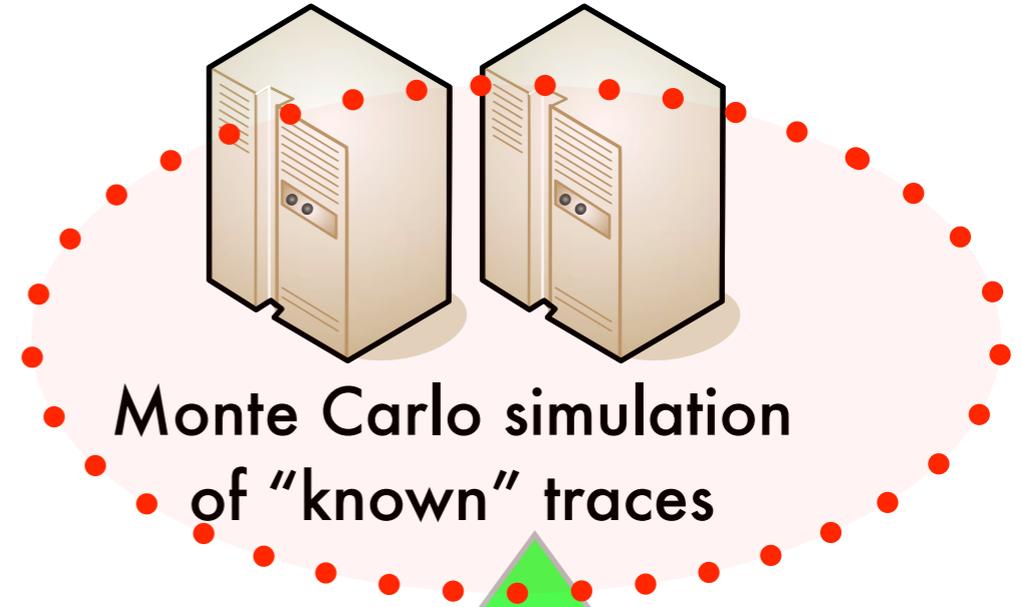
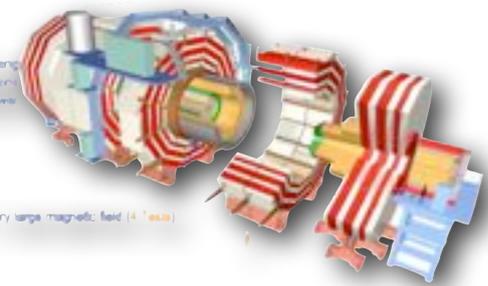
- CMS: Compact Muon Solenoid at CERN
 - 3500 scientists, 180 Universities and Research Labs (40 countries)
 - CMS is like a ~75 MegaPixels Digital Camera. 40M “photos” /s
Selection of 300 ‘photos’ /s ~450 MB /s from the detector are ~PBs of data /year
 - CERN has (of course) its well established data flow and infrastructure, however ...



Tier-0 <i>(the accelerator centre)</i>	7 Tier-1s <i>("online" to the DAQ)</i>	~50 Tier-2s <i>in ~20 countries</i>
Data acquisition & initial processing Long-term mass data storage CMS CERN Analysis Facility <i>(latency-critical data processing, high priority analysis)</i> Distribution of data → Tier-1 centres	High availability centers Custodial mass storage of share of data Data reconstruction and reprocessing Data skimming Distribute analysis data → Tier-2s	End-user physics analyses Detector Studies MonteCarlo Simulation → Tier-1

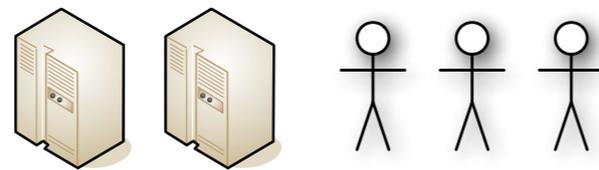
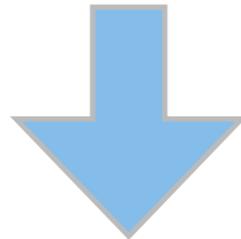


Innovative Methods for Particle Colliders at the Terascale (2012-2015, oversimplified vision)

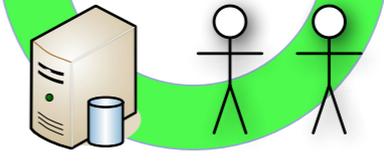
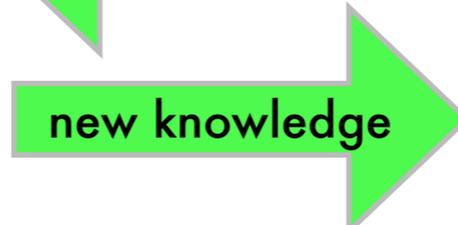
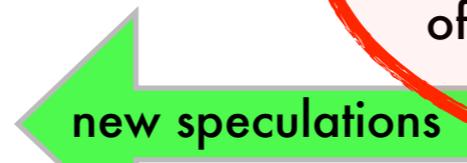
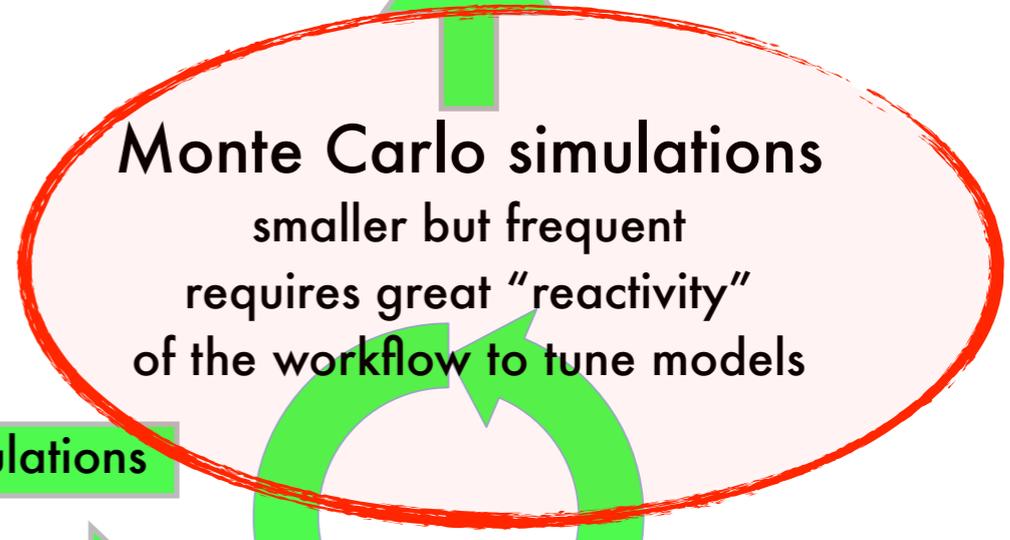


Preliminary results on channel $H \rightarrow ZZ \rightarrow 4l$ shown at Higg's boson claim Jul 2012 (Amapane et al.)

Filtered data to be analysed (much smaller)



Particle physicists



Theoretic physicists

Formalising the cell cycle switch

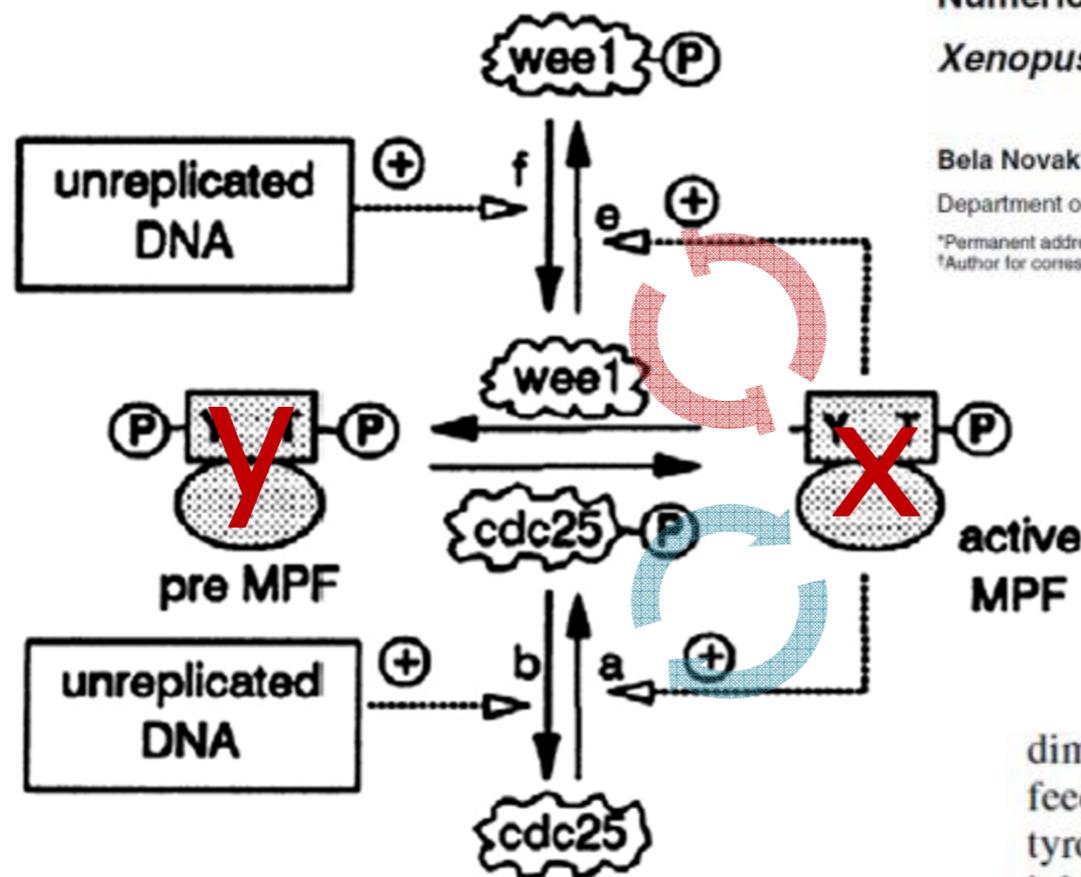
Journal of Cell Science 106, 1153-1168 (1993)
 Printed in Great Britain © The Company of Biologists Limited 1993

Numerical analysis of a comprehensive model of M-phase control in *Xenopus* oocyte extracts and intact embryos

Bela Novak* and John J. Tyson†

Department of Biology, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24060-0406, USA

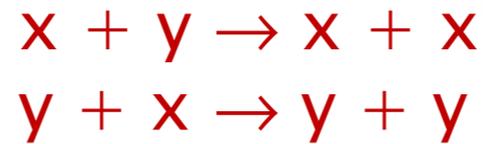
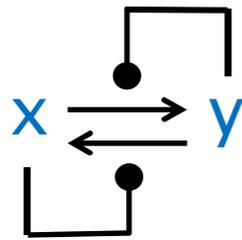
*Permanent address: Department of Agricultural Chemical Technology, Technical University of Budapest, 1521 Budapest Gellert Ter 4, Hungary
 †Author for correspondence



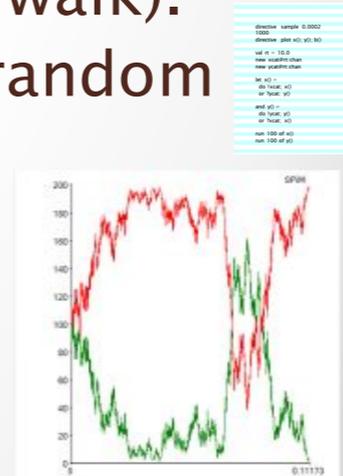
dimers is left off the diagram to keep it simple.) (B) Positive feedback loops. Active MPF stimulates its own production from tyrosine-phosphorylated dimers by activating Cdc25 and inhibiting Wee1. We suspect that these signals are indirect, but intermediary enzymes are unknown and we ignore them in this paper. The signals from active MPF to Wee1 and Cdc25 generate an autocatalytic instability in the control system. We indicate also an 'external' signal from unreplicated DNA to Wee1 and Cdc25, which can be used to control the efficacy of the positive feedback loops. The letters a, b, e and f are used to label the rate constants for these reactions in Fig. 2. (C) Negative feedback loop. Active

Direct competition: unstable switch

- x catalyzes the transformation of y into x
- y catalyzes the transformation of x into y



- This system is bistable, but
 - Convergence to a stable state is slow (a random walk).
 - *Any* perturbation of a stable state can initiate a random walk to the other stable state.
 - With 100 molecules of x and y, convergence is quick, but with 10000 molecules, even at the same concentration (adjusting the rate) you will wait for a long time.

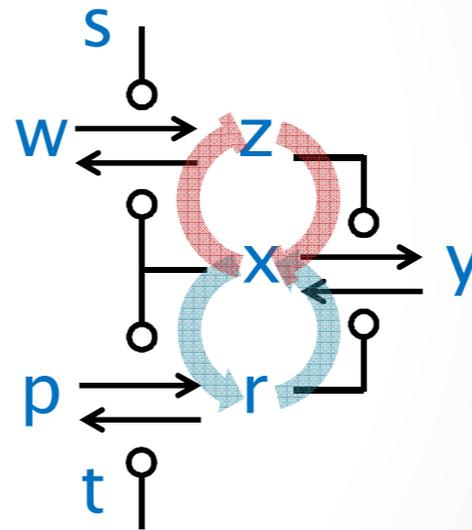
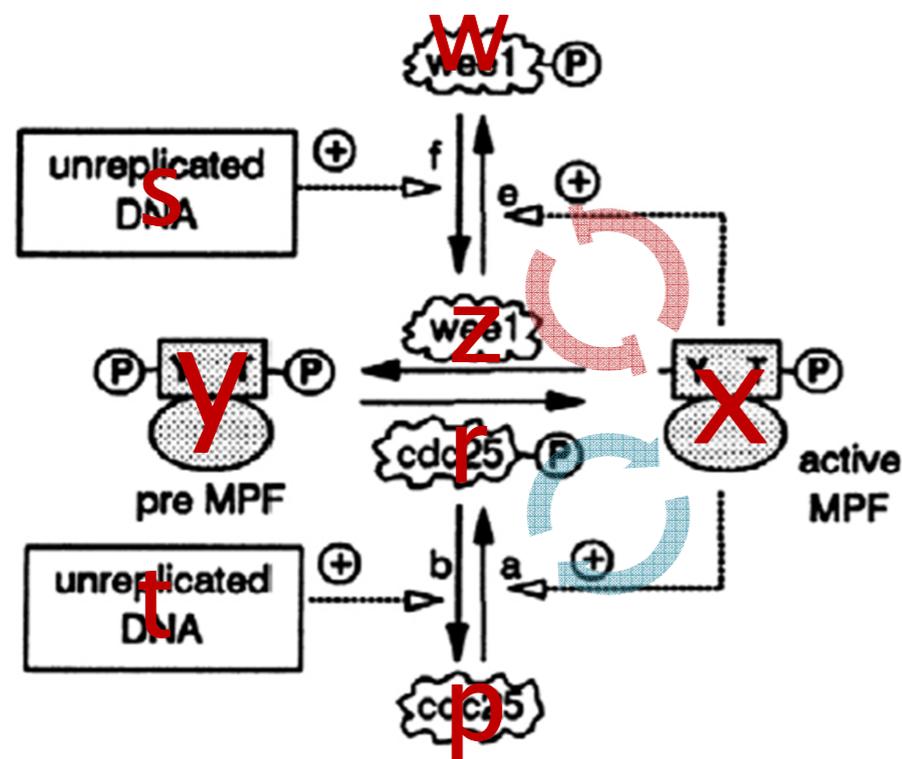


Courtesy of Luca Cardelli
 On Switches and Oscillators Program
 Equivalence in Biology?
<http://luccardelli.name>

CWC syntax



... after a number of transformations: a stable switch faithfully modelling cell switch



Courtesy of Luca Cardelli

On Switches and Oscillators Program
Equivalence in Biology?

<http://lucacardelli.name>

CWC
syntax

$$T : a c \xrightarrow{10} c b$$

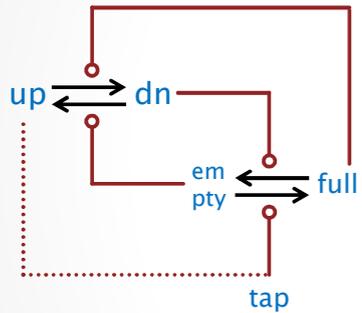
$$T : c a \xrightarrow{10} a b$$

$$T : b a \xrightarrow{10} a a$$

$$T : b c \xrightarrow{10} c c$$

The Shishi Odoshi

- A Japanese scarecrow (scare-deer)
 - Used by Bela Novak to illustrate the cell cycle switch.



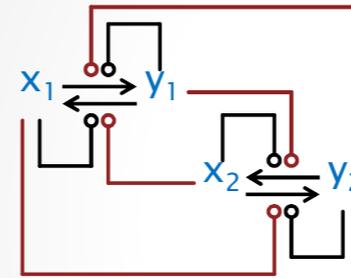
<http://www.youtube.com/watch?v=VbvecTiftcE&NR=1&feature=fwvp>

empty + tap \rightarrow tap + full
 up + full \rightarrow full + dn
 full + dn \rightarrow dn + empty
 dn + empty \rightarrow empty + up

To make it into a full trammel (dotted line), we could make the up position mechanically open the tap (i.e. take up = tap)

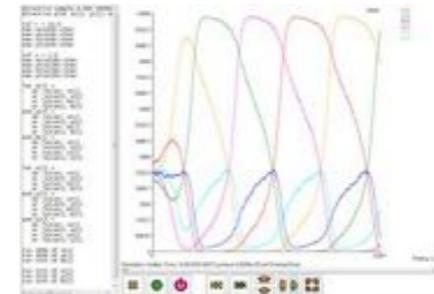
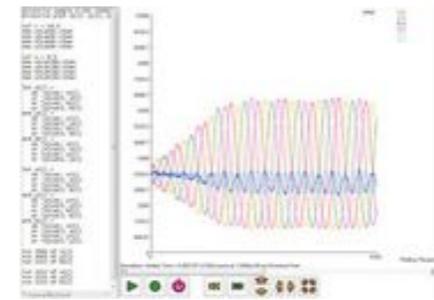
The 2AM Limit-Cycle Oscillator

- Two AM switches in a Trammel pattern



The red reactions need to be slower (even slightly) than the black reactions, but otherwise the oscillation is robust. Oscillation stops at 10 vs. 10 and 1 vs. 10. Here the rates are 8 vs 10.0 top, and 2 vs 10, bottom.

(Simple limit-cycle oscillators in the literature have very critical rate ranges.)



```
directive sample 0.001 10000
directive plot x10; y10; b10; x20;
y20; b20

val r = 10.0
new x1cat@r:chan
new y1cat@r:chan
new x2cat@r:chan
new y2cat@r:chan

val s = 8.0
new x1cat2@s:chan
new y1cat2@s:chan
new x2cat1@s:chan
new y2cat1@s:chan

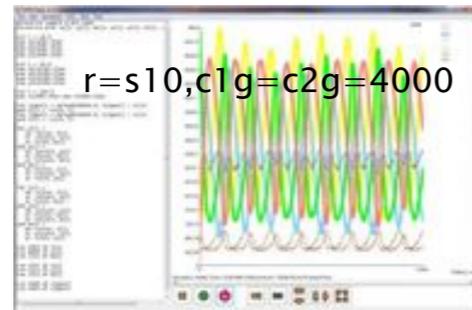
let x10 =
do lx1cat: x10
or lx2cat2: x10
or y1cat: b10
or y2cat1: b10
and y10 =
do y1cat: y10
or y1cat2: y10
or x1cat: b10
or x2cat1: b10
and b10 =
do x1cat: x10
or x2cat1: x10
or y1cat: y10
or y2cat1: y10

let x20 =
do lx2cat: x20
or lx2cat1: x20
or y2cat: b20
or x1cat2: b20
and y20 =
do y2cat: y20
or y2cat1: y20
or x2cat: b20
or y1cat2: b20
and b20 =
do x2cat: x20
or y1cat2: x20
or x2cat1: y20
or x1cat2: y20

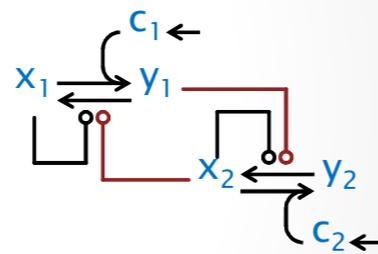
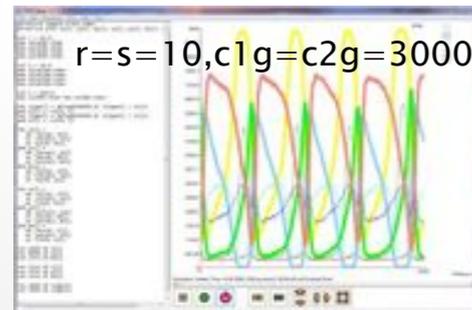
run 3666 of x10
run 3000 of y10
run 3333 of b10
run 3333 of x20
run 3333 of y20
run 3333 of b20
```

Influx Oscillators

- Similar but:
 - The two-input switches are replaced by one-input switches which are reset by constant influxes.



$r=s=10, c1g=c2g=3000$

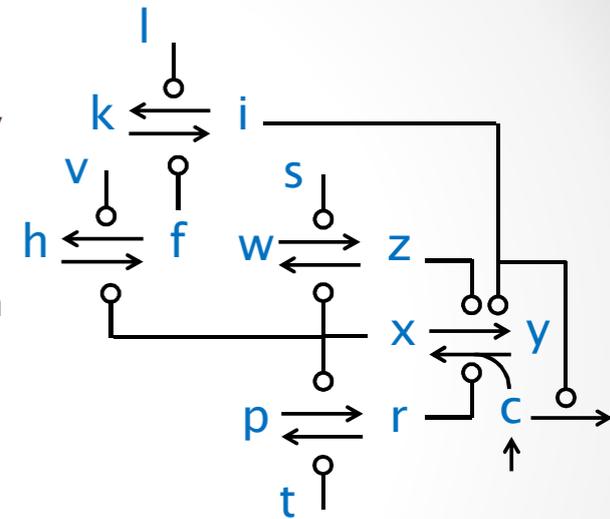


Works best with $s=r$.

Needs constant influx of $c1, c2$

Novak-Tyson Oscillator

- First switch
 - Is the 'transformed' AM switch in one-input configuration (driven by constant influx of cyclin).
- Second switch
 - Is a simple two-stage switch working as a delay (the first switch is so good in terms of hysteresis that the second switch is not very critical for oscillation).
 - It can be replaced by a one-stage switch (Ferrell's cell cycle oscillator) but oscillation is a bit harder to obtain.
- Connection
 - Single links, as in the influx oscillator.

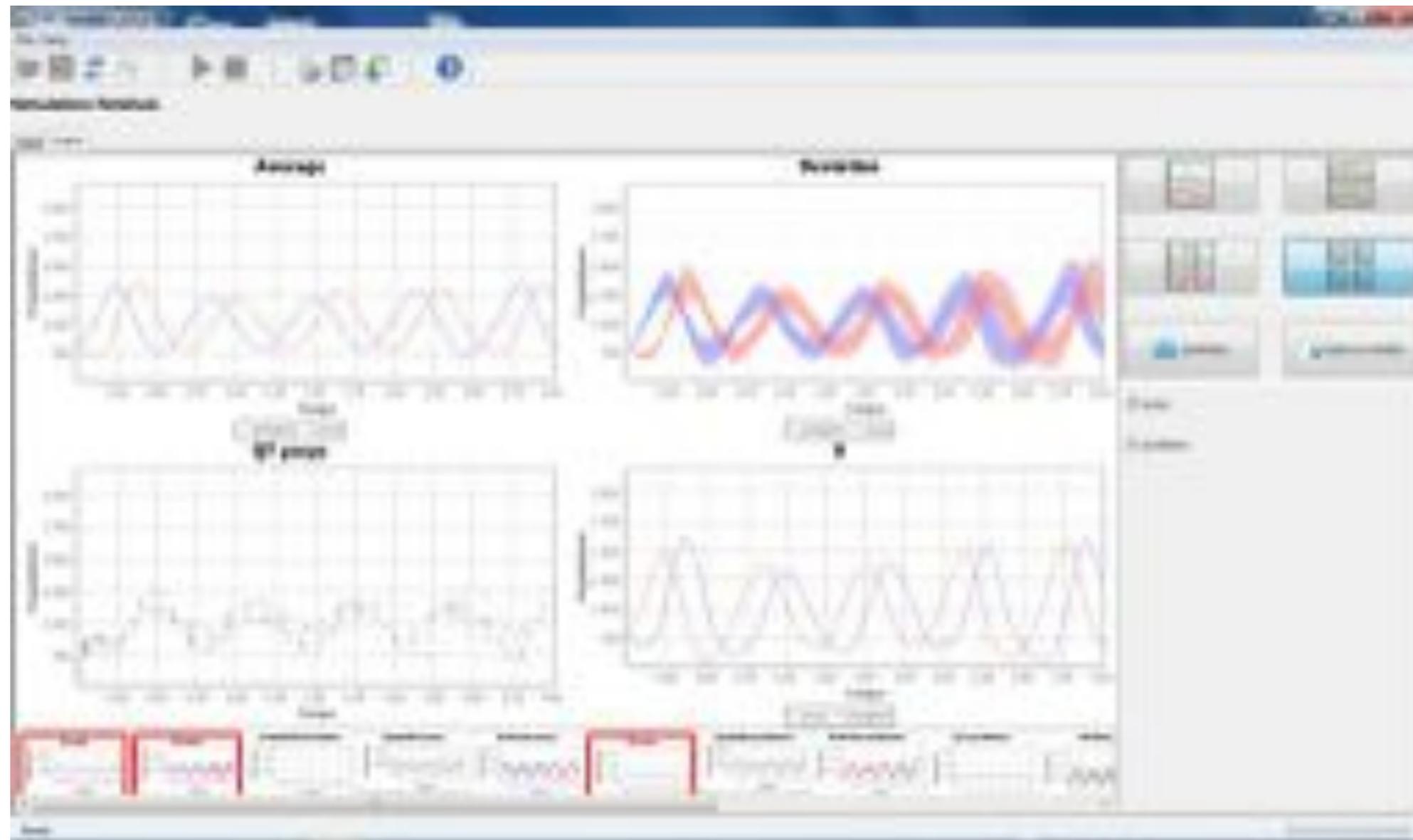


Journal of Cell Science 106, 1153-1161 (1993)
 Printed in Great Britain © The Company of Biologists Limited 1993

Numerical analysis of a comprehensive model of M-phase control in *Xenopus* oocyte extracts and intact embryos

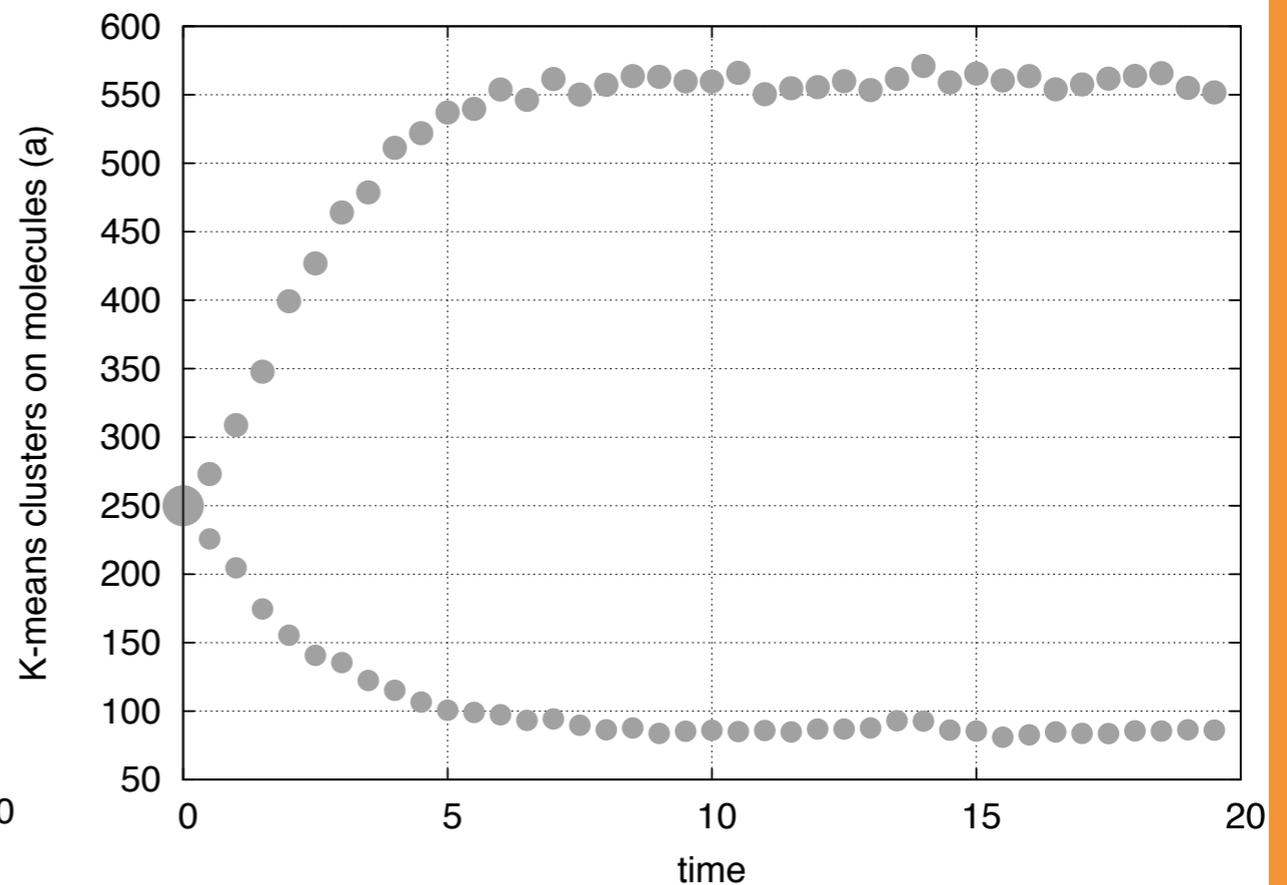
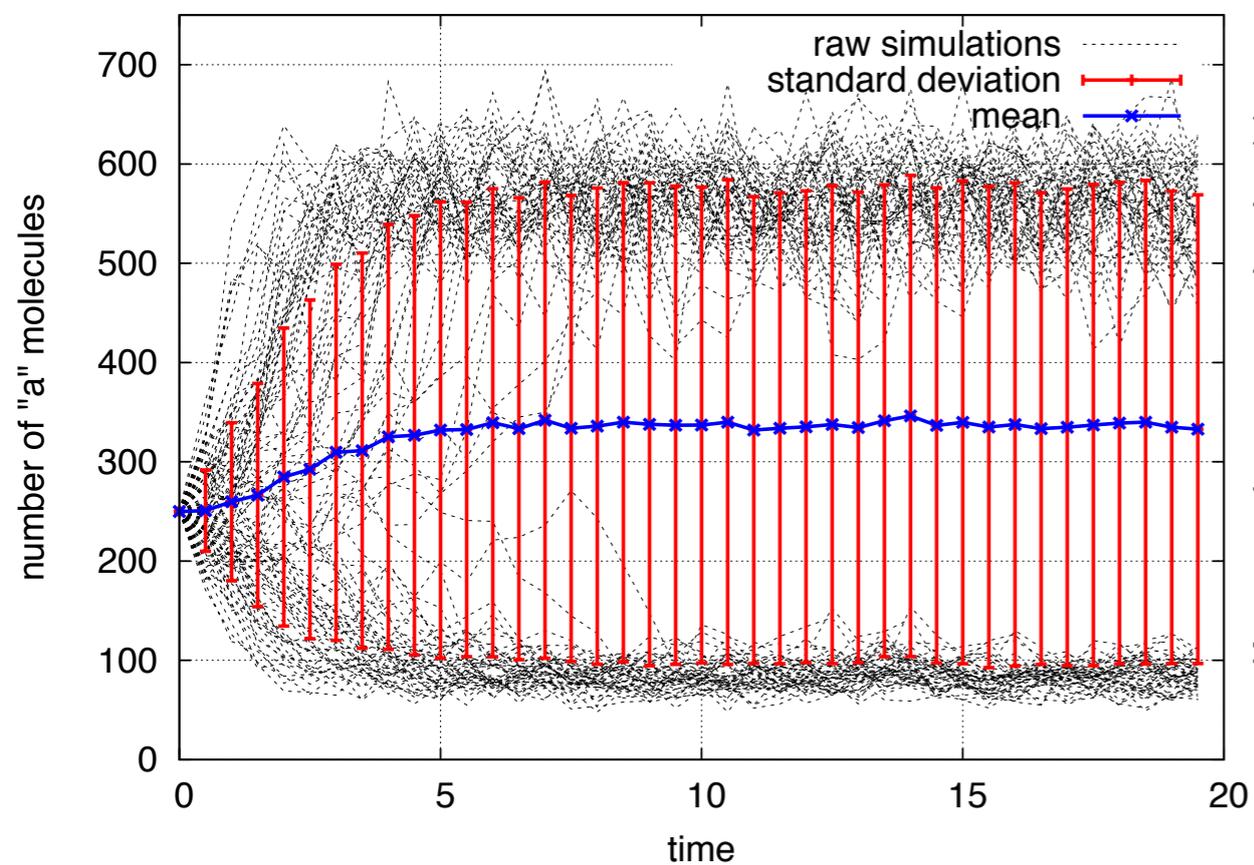
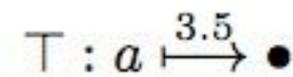
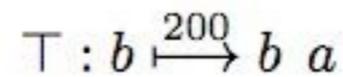
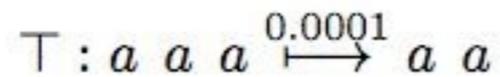
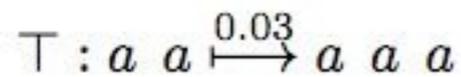
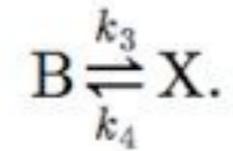
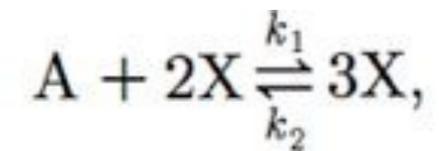
Bela Novak* and John J. Tyson†
 Department of Biology, Virginia Polytechnic Institute and State University, Blacksburg, Virginia 24060-0406, USA

- Demo: unstable switch

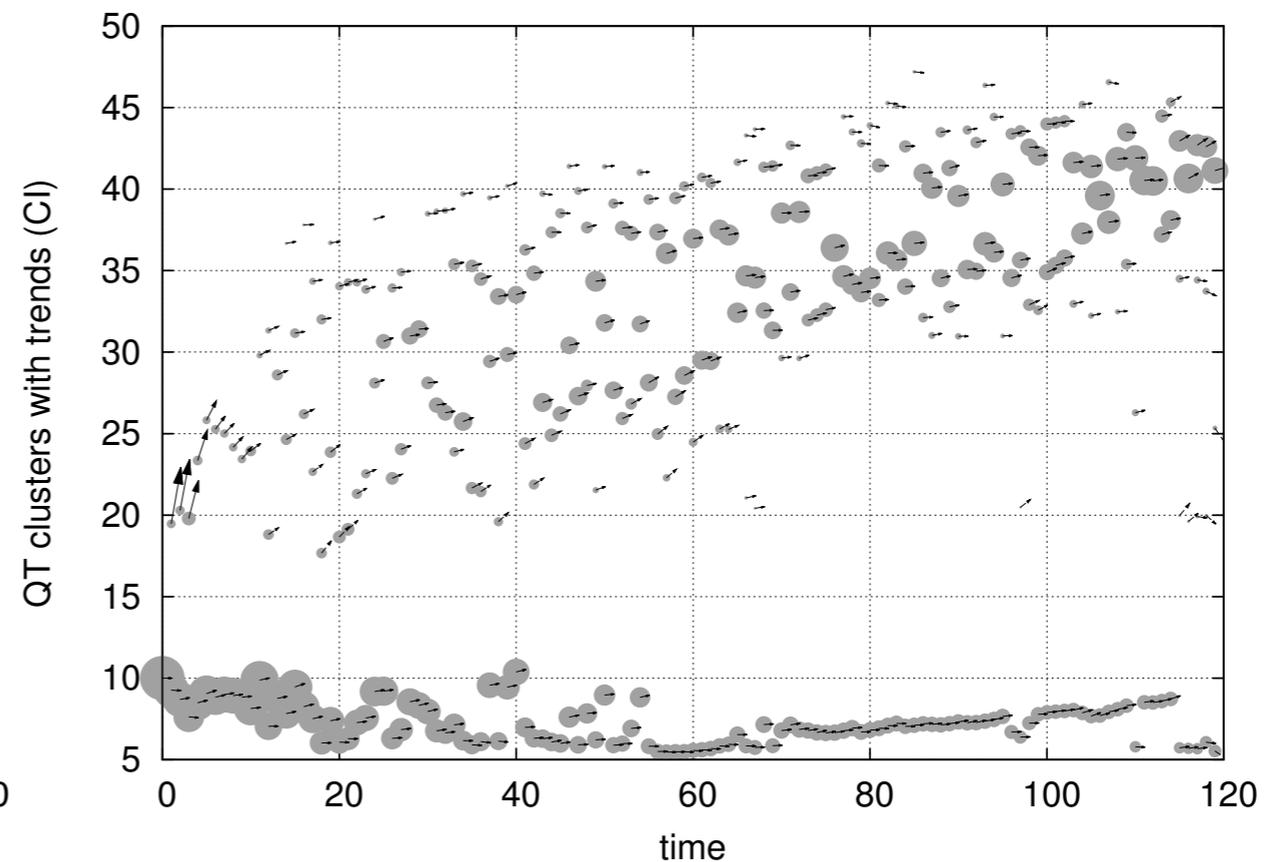
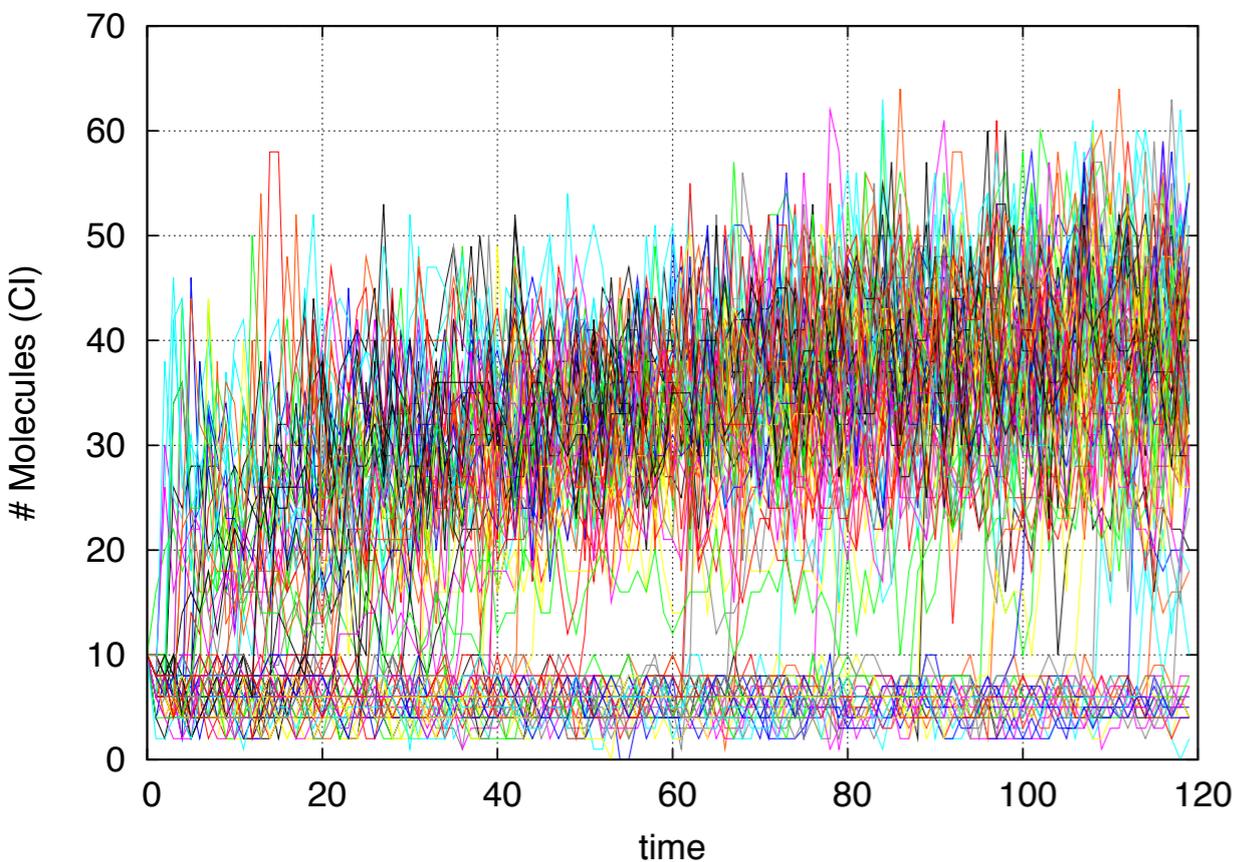
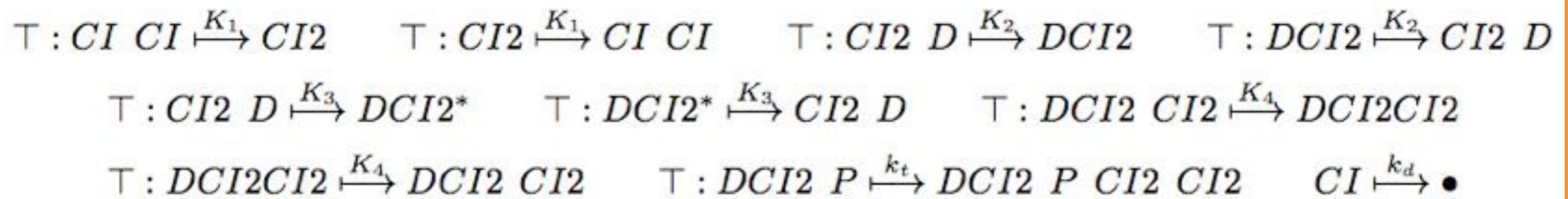
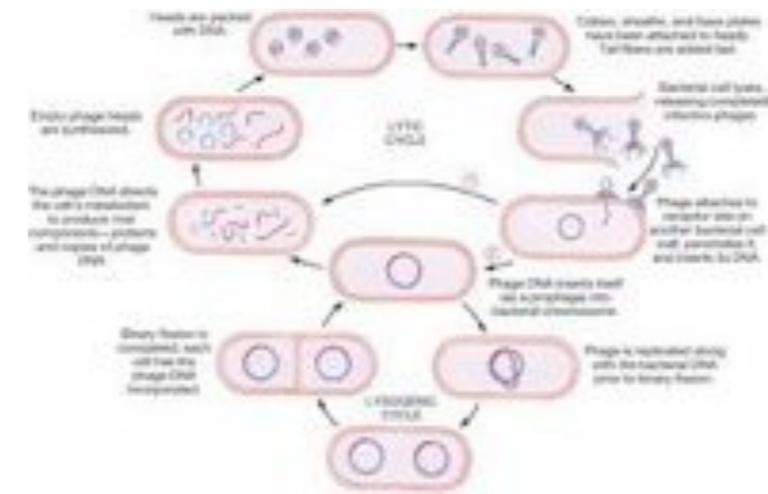


Schlögl model

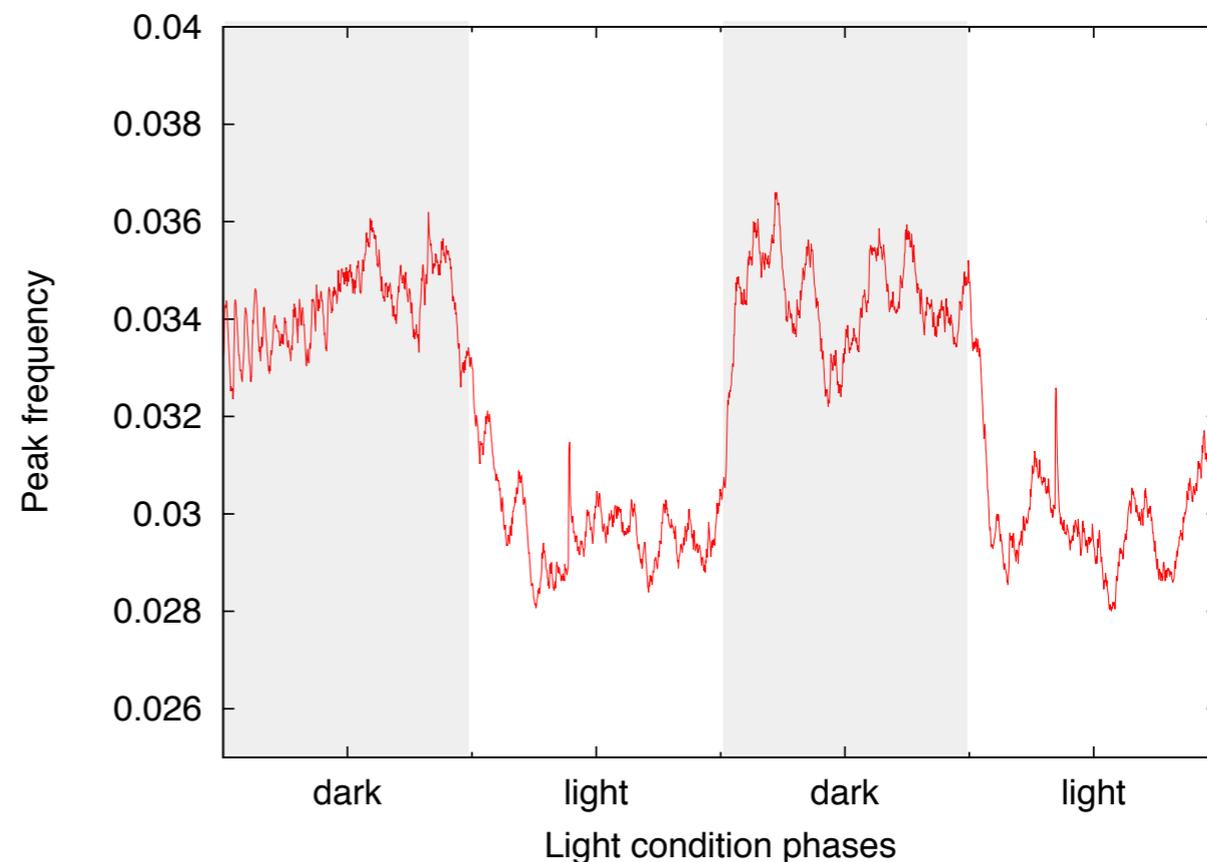
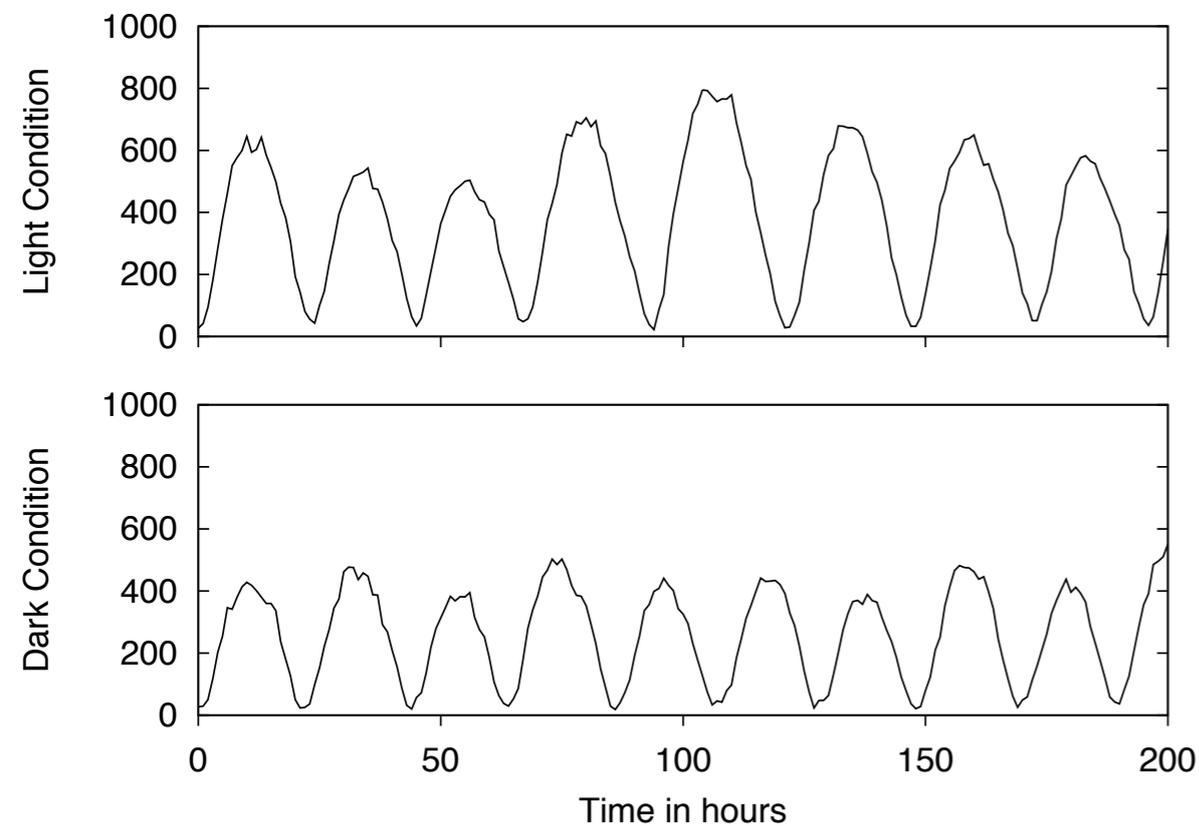
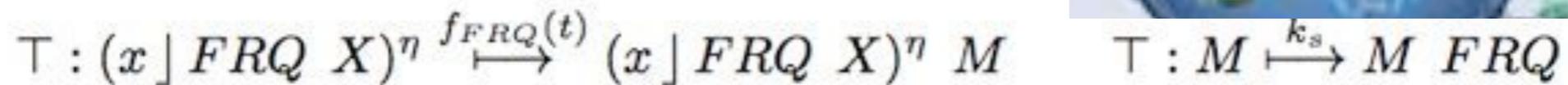
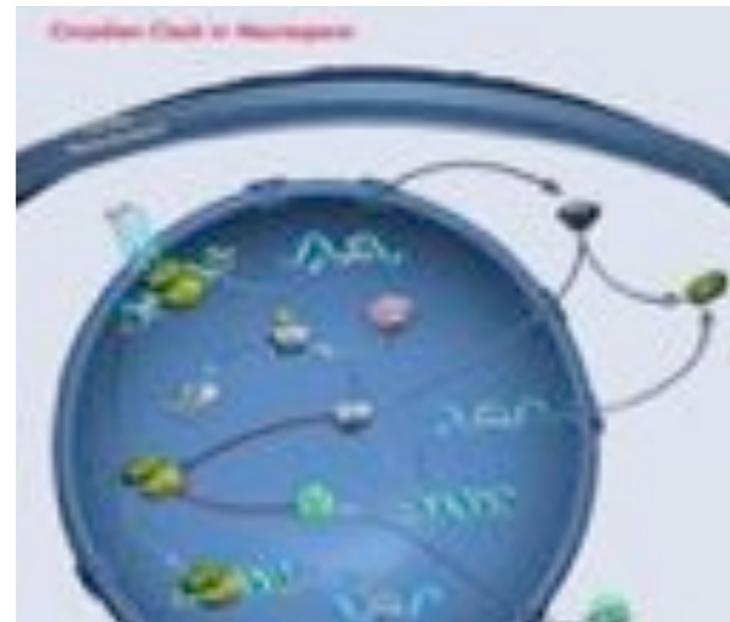
autocatalytic, trimolecular reaction scheme (bistable)



Bacteriophage λ life cycle integration of a strand of DNA in the molecule of E. coli DNA (multi-stable)



Transcriptional regulation in Neurospora (circadian clock period detection)



Conclusions

- Talk focused on programming model
 - Many important “in a cloud” ignored in the talk: middleware, faults, ...
 - Data movement & high-level are key features
 - helps the mapping of features to platforms and performance portability
 - ease the design
 - MapReduce is an instance, should be not the only one
- Formal biology at embryonal stage
 - similar data & computation problems in analysis “in vitro” experiments
 - increasing interest from industrial and “core” bio scientist for parallel computing

Acknowledgements

- FastFlow
 - Massimo Torquati (Pisa), Marco Danelutto (Pisa), Peter Kilpatrick (Belfast), Massimiliano Meneghin (IBM Research Dublin)
- Case studies, biological models
 - Luca Cardelli (Microsoft), Pietro Liò (Cambridge), Andrea Bracciali (Stirling), Cristina Calcagno (Torino)
- CWC simulation implementation
 - Maurizio Drocco (Torino), Fabio Tordini (Torino)
- CWC language design
 - Maurizio Drocco, Eva Sciacca, Salvatore Spinella, Mario Coppo, Angelo Troina, Ferruccio Damiani (Torino)
- Graphical User Interface
 - ETICA company
- Infiniband cluster and computing facilities
 - Mellanox, HPC Advisory Council

