

Going Through Rough Times:

Parallel Discrete Event Simulations (PDES)
--- A Physicist's Perspective

Mark A. Novotny

Dept. of Physics and Astronomy

HPC² Center for Computational Sciences

Funded in part the NSF

(DMR9871455) $\xrightarrow{\text{ITR}}$ DMR0113049 $\xrightarrow{\text{ITR}}$ DMR0426488



Collaborators

Alice Kolakowska (Purdue U, Calumet, former post-doc)

Gyorgy Korniss (Rensselaer Polytechnic Institute)

Zoltan Rácz (Eötvös University, Hungary)

Per Arne Rikvold (Florida State University)

Lev Shchur (Scientific Center in Chernogolovka, Russia)

Zoltan Toroczka (Notre Dame University)

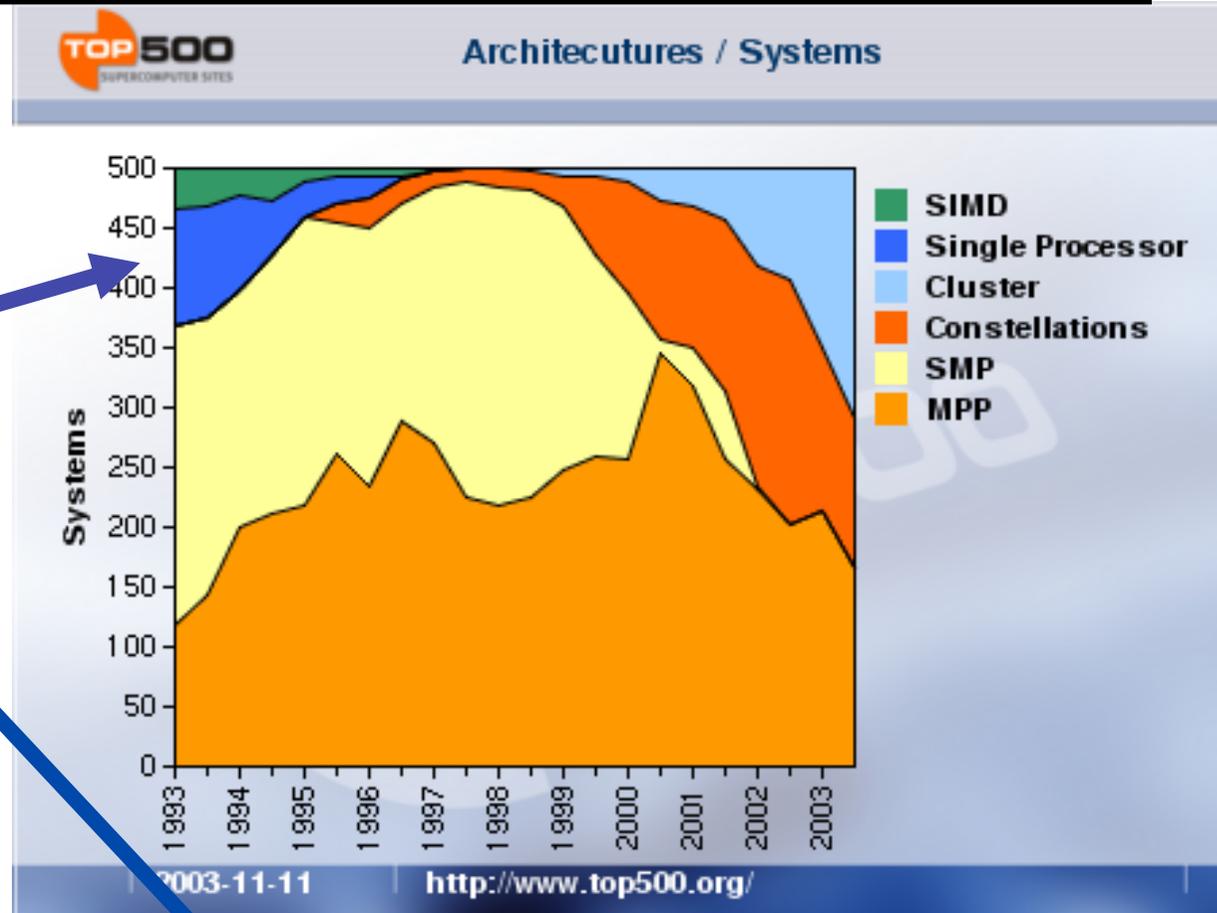
Students

H. Guclu, B. Kozma (Rensselaer Polytechnic Inst.)

Terrance Dubrues, Poonam Verma (Mississippi State University)

Daniel Brown, Melissa Cook, Sara Gill, Tori Norwood,
Amanda Novotny (summer High School students)

Problem: no signal faster than speed of light

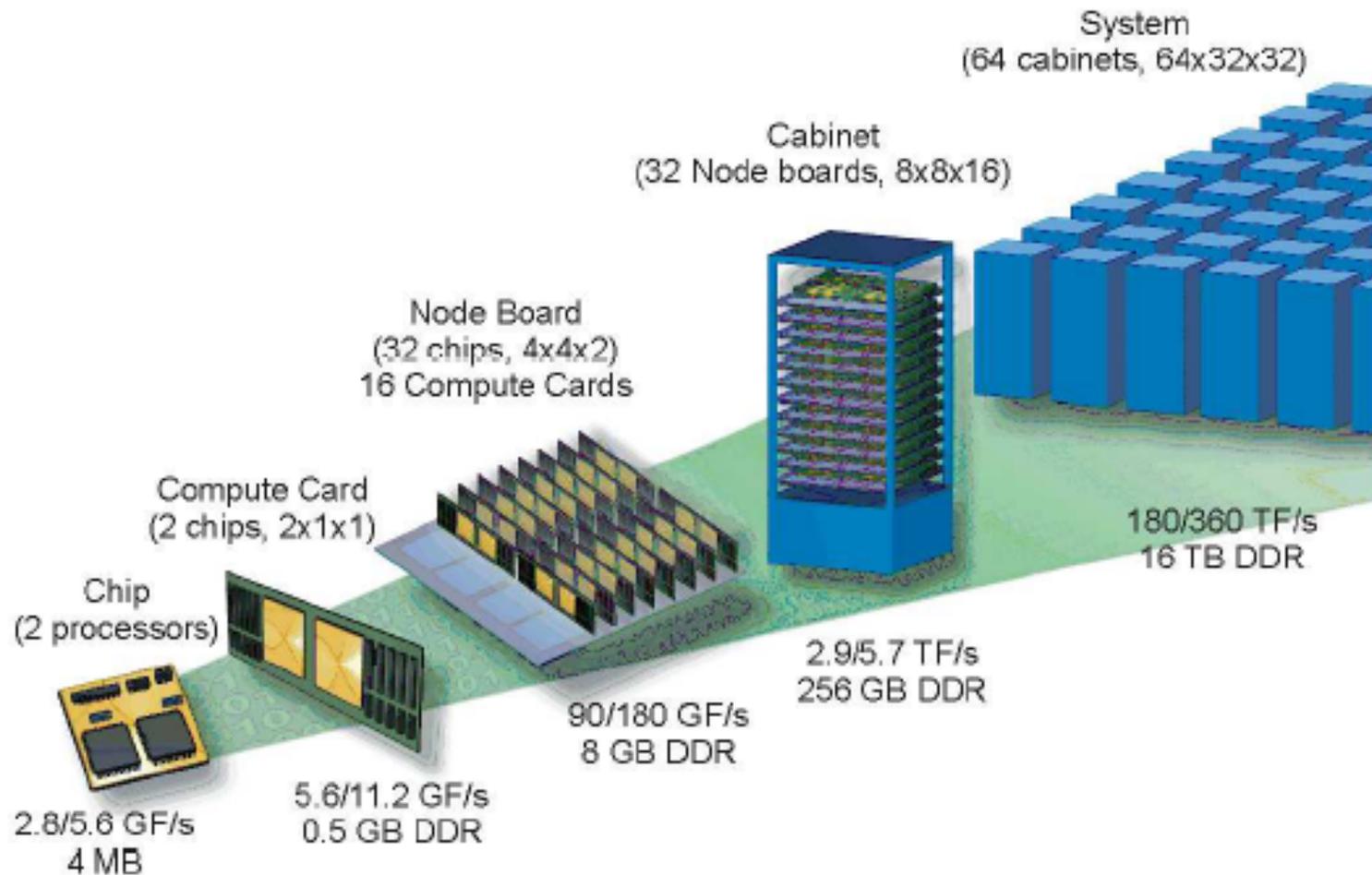


Since 1996 (at least)

computing on supercomputers *requires* parallel computing

Prediction: in 10 years almost all PCs will have >16 PEs
(Processing Elements)

BlueGene/L → BlueGene/Q



$$2 \times 2 \times 32 \times 32 \times 64 = 131,072 \text{ PEs}$$

June 2012 Top500.org

Sequoia computer
BlueGene/Q IBM
US DOE NNSA/LLNL
1,572,864 cores
 $R_{\text{peak}} = 15.7$ petaFLOP
 $R_{\text{max}} = 20.1$ petaFLOP



**Into the Wide
Blue Yonder
with BlueGene/L**

0 with less than 1024 Pes
1 Top-500 with 1024 to 2048 PEs
12 Top-500 with $\geq 128,000$ PEs

Perfect Scalability

N_{PE} = # Processing Elements

Non-trivial parallelization

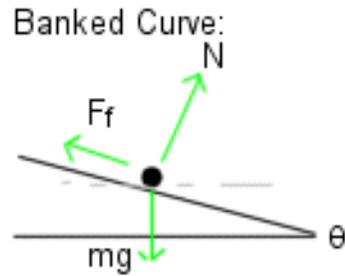
Amount of work $O(N_{PE})$

- For large N_{PE} , utilization independent of N_{PE}
- For large N_{PE} , PE memory independent of N_{PE}
- Number of interconnects $O(1)$ per PE

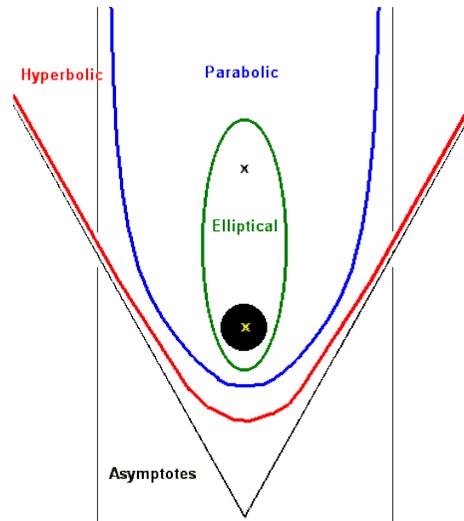


How Computational Physicists Count

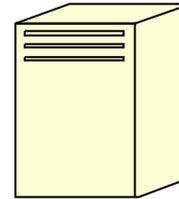
1 body



2 body



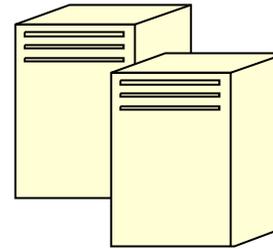
Too Many bodies
use
Statistical Mechanics



1 PE (Processor Element)

=

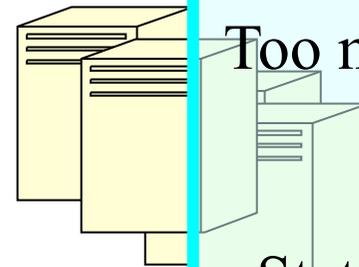
1 grad student



2 PEs

=

2 grad students



Too Many PEs

Too many grad students

use

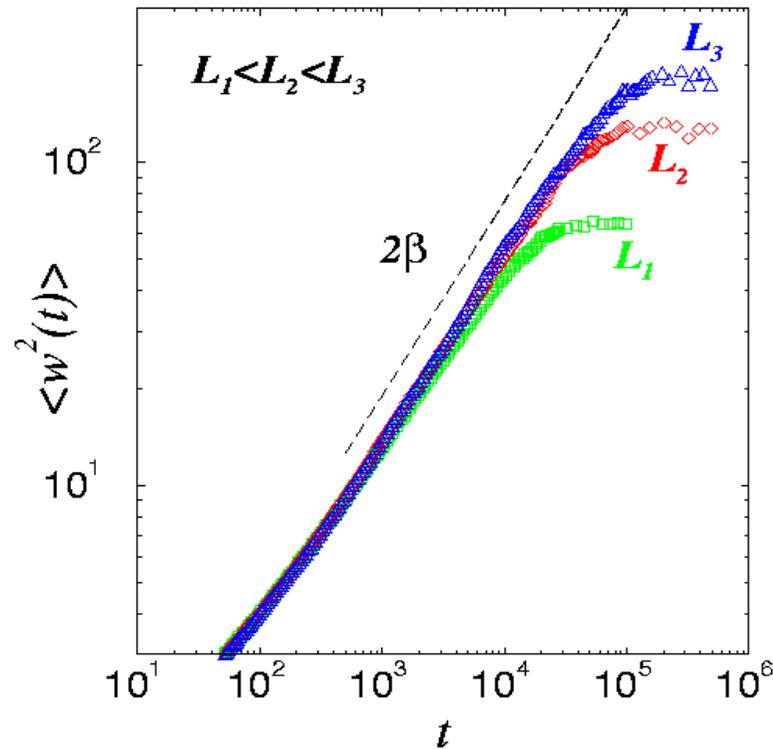
Statistical Mechanics

Complicated Behavior & Informatics
from Non-equilibrium Surface Growth Models

Motivation for PDES model

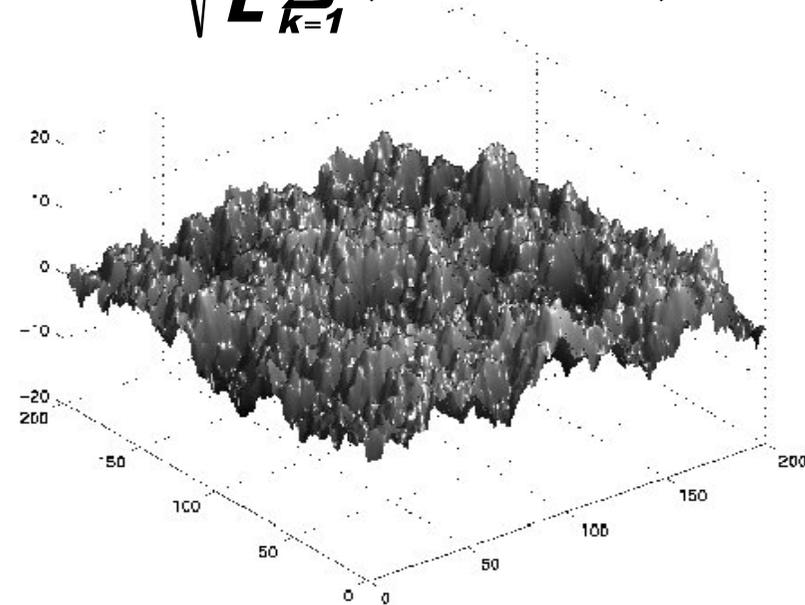
Parallel computing

Non-equilibrium surface growth



$$\langle w^2(t) \rangle \sim \begin{cases} t^{2\beta}, & \text{if } t \ll t_x \\ L^{2\alpha}, & \text{if } t \gg t_x \end{cases}$$

$$w(\mathbf{t}) = \sqrt{\frac{1}{L} \sum_{k=1}^L (\tau_k(\mathbf{t}) - \bar{\tau}(\mathbf{t}))^2}$$



Dynamic scaling:

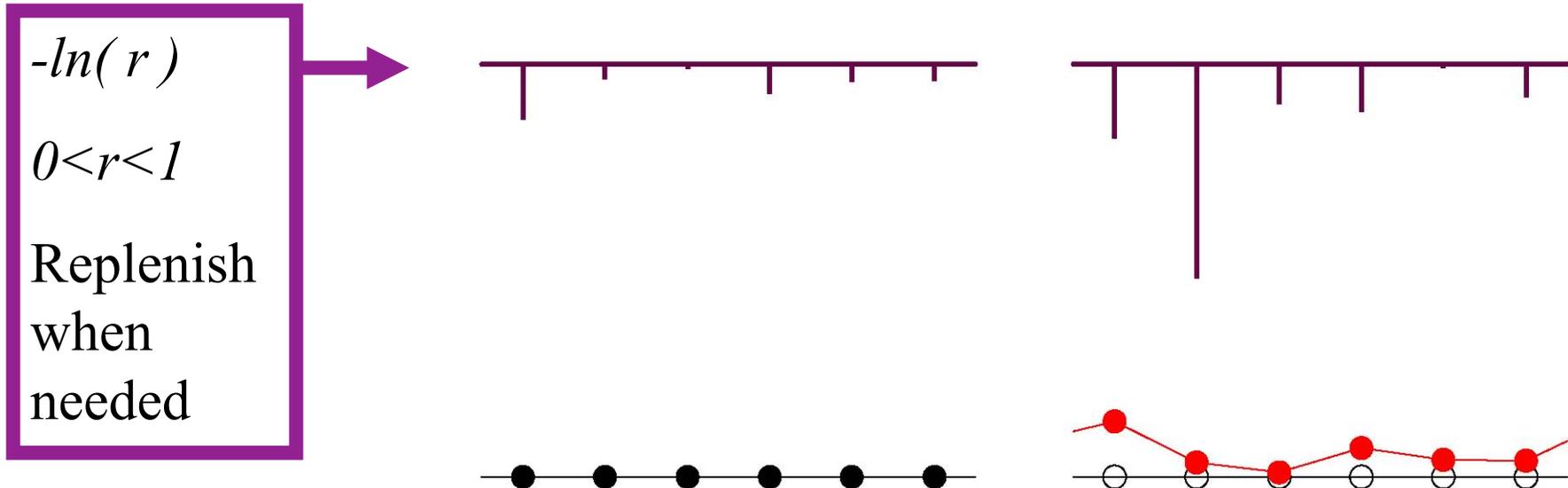
$$\alpha = \beta z$$

β growth exponent

z dynamic exponent

α roughness exponent

Non-equilibrium surface growth model: PDES model



Start with flat interface (*in d dimensions*)

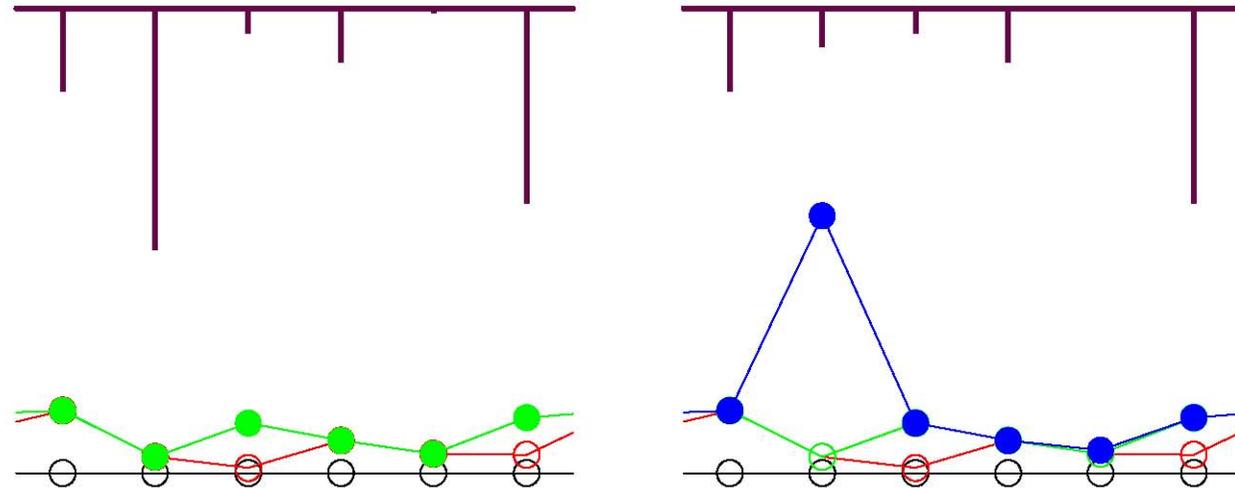
In first step, all 'drops' fall

PDES *model*

$$-\ln(r)$$

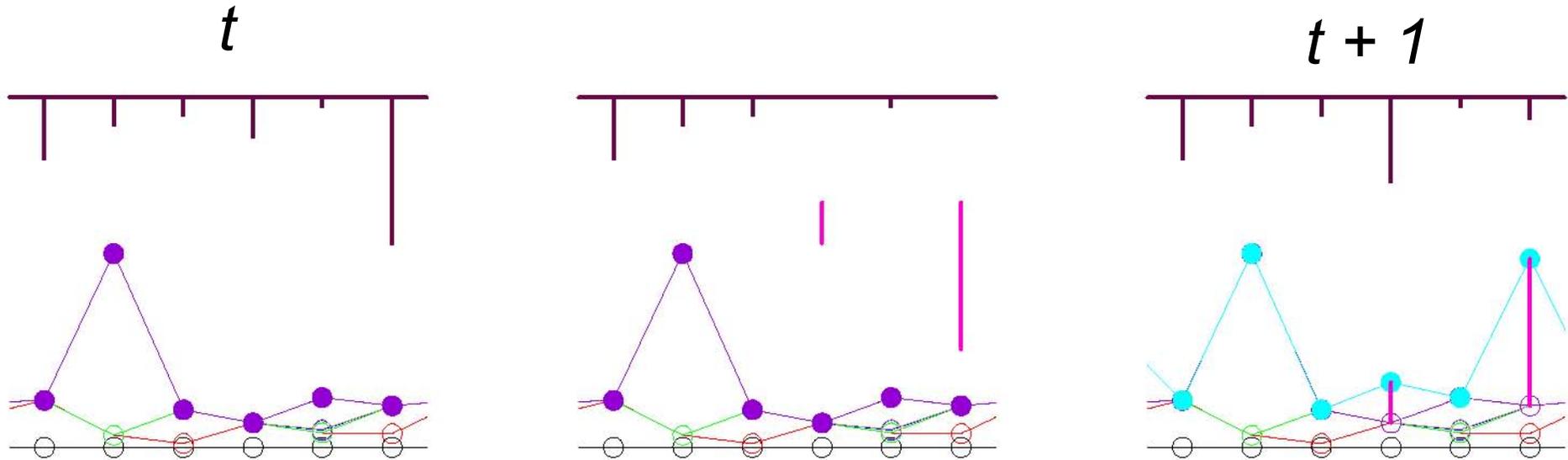
$$0 < r < 1$$

Replenish
when
needed



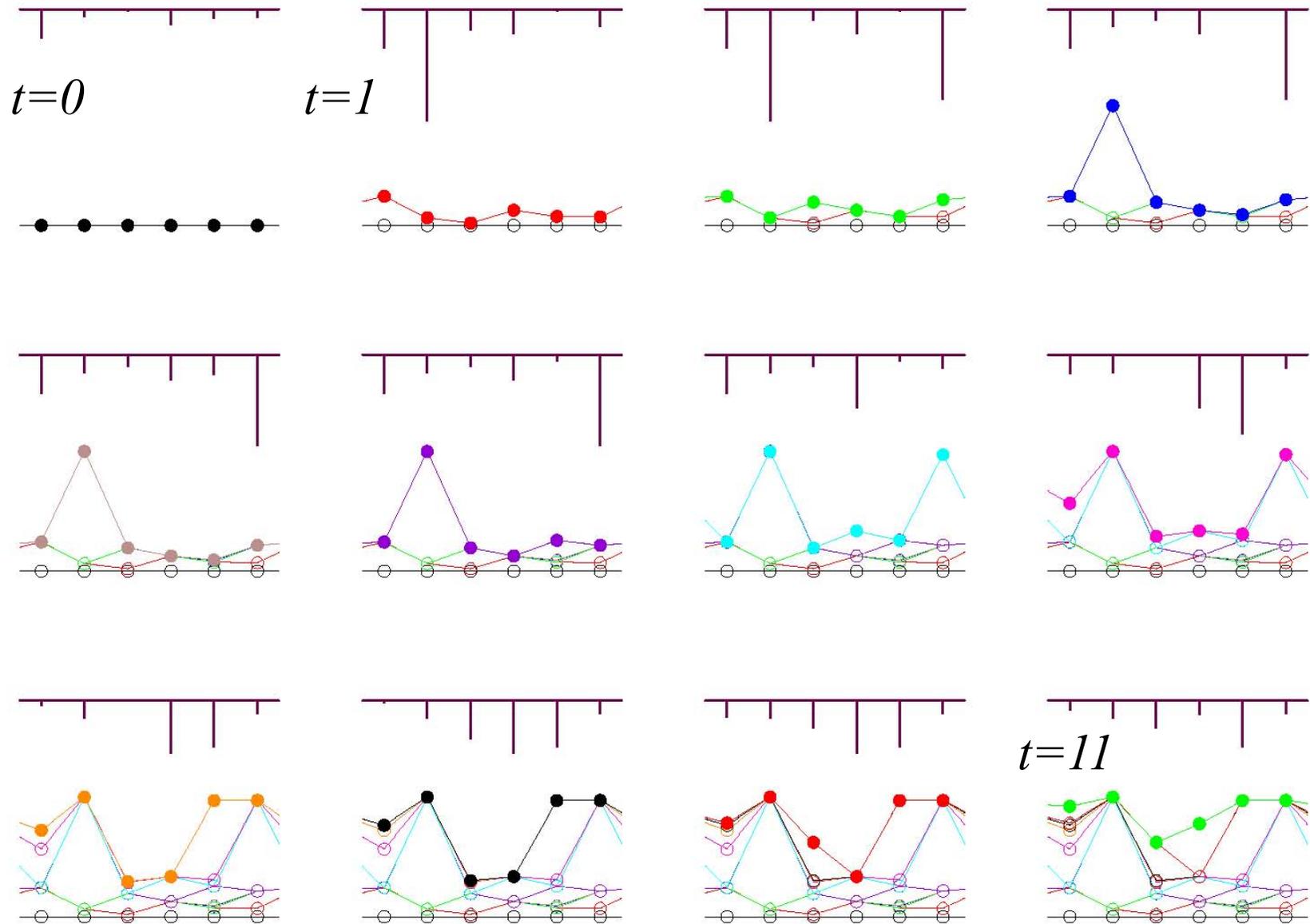
For each step, all 'drops' fall ***ONLY*** if the surface underneath is at a local minimum

PDES *model*



Note: at each step t all ‘drops’ fall
at the same time

PDES *model*



Discrete Event Simulations



- DES (Discrete Event Simulations)
 - * State changes are discontinuous
 - * Times of state changes are random

PDES

Parallel Discrete Event Simulations

PDES Technology Implications

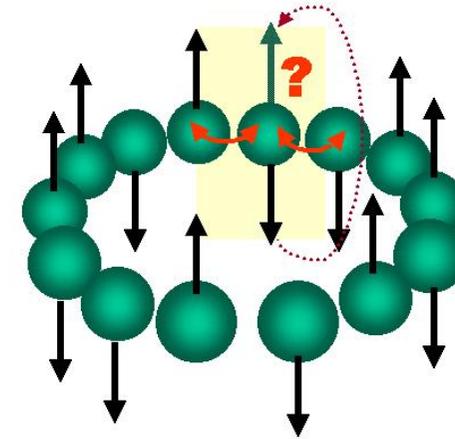
- **All today's largest computers are massively parallel computers**
- **Must make good use of parallelization in programs for efficiency**
- **Parallel Discrete Event Simulations (PDES)**
 - **Used in military simulations and training ('what-if' scenarios)**
 - **Used in homeland security simulations and training**
 - **Used in modeling of factory deliveries**
 - **Used in modeling temporal drug concentrations in patient models**
 - **Used in simulating materials and materials failure**
 - **Used in modeling switching in cellular and wireless networks**
 - **Used in ecological modeling**
 - **Used in modeling epidemiological models**
 - **Used in electric power grid simulations**

Information-Driven Systems

Example:
Dynamic Monte Carlo of Ising spins
with nearest-neighbor interactions

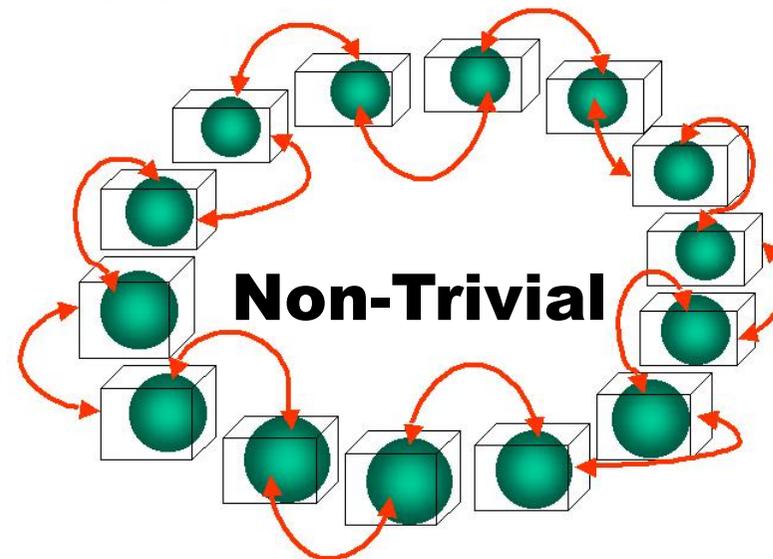
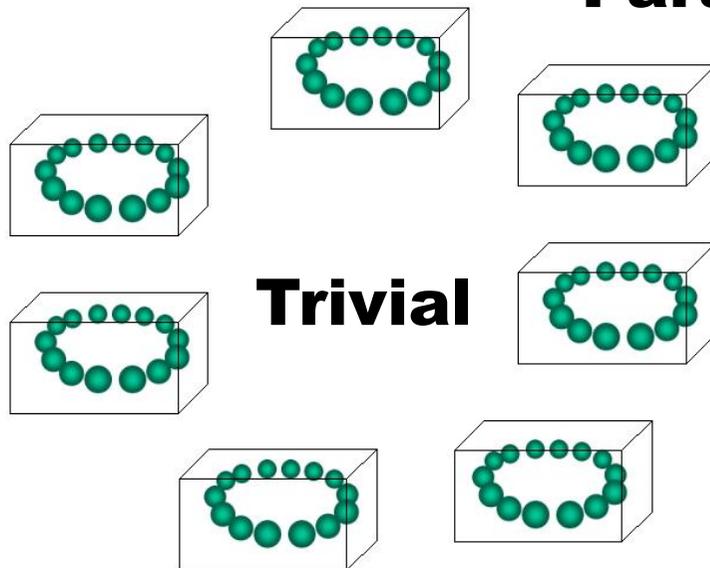
Randomly pick a spin

Decide if spin will be flipped



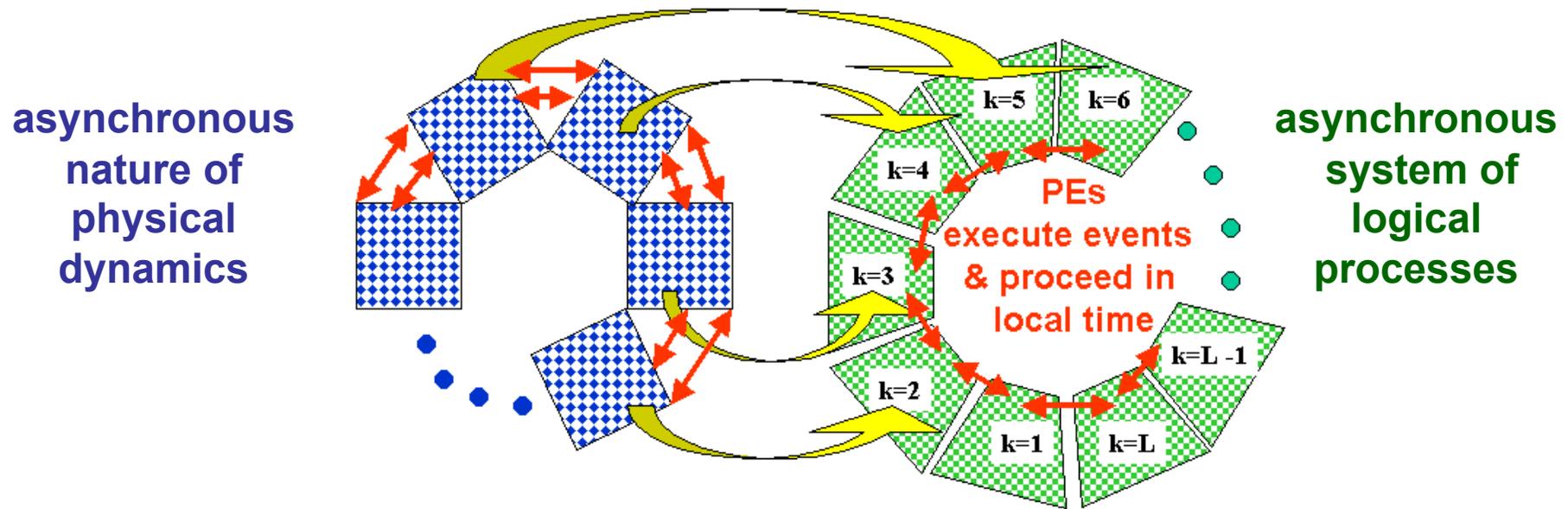
Dynamic Monte Carlo simulations

Parallelization





Physical processes and logical processes



Physical System
spatially extended system of NL spins, arranged on a lattice

Computing System
 L PEs: each carries N lattice sites, N_b of which are border sites

Physical Events/Processes
random spin flipping

Logical Events/Processes
each PE manages the state of the assigned subsystem.

discrete event: the spin flip

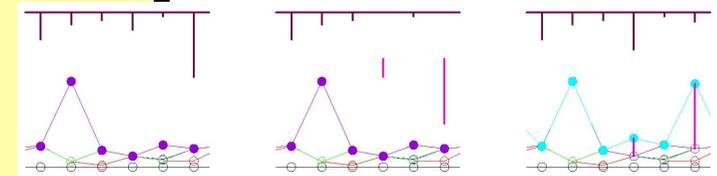
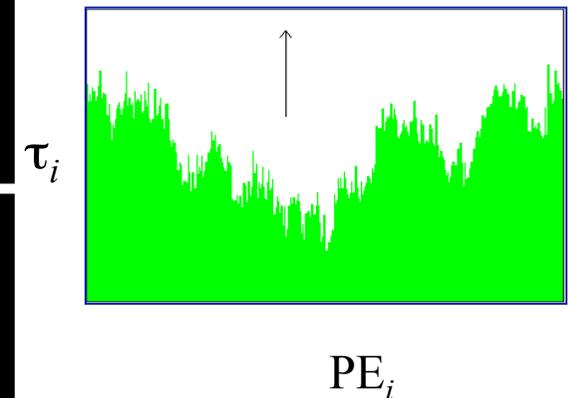
discrete event: the state update

Parallel discrete-event simulation

for spatially decomposable **asynchronous cellular automata**

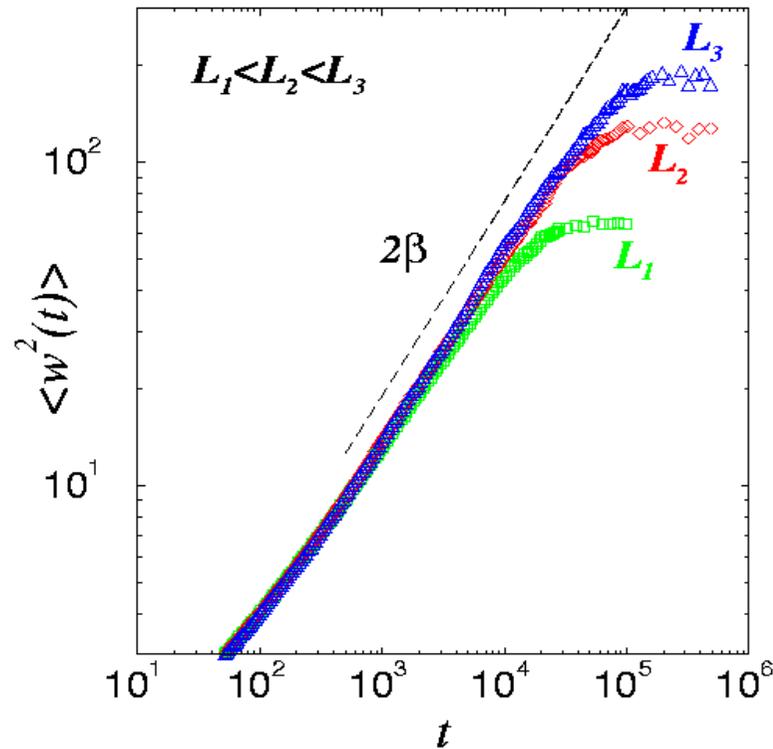
- **Spatial decomposition** on lattice/grid
(for systems with **short-range interactions**
only **local synchronization** between subsystems)
- Changes/updates: independent Poisson arrivals

- ❖ Each subsystem/block of sites, carried by a processing element (PE) must have its own **local simulated time, $\{\tau_i\}$** (“virtual time”)
- ❖ **Synchronization** scheme
- ❖ PEs must concurrently advance their own Poisson streams, **without violating causality**



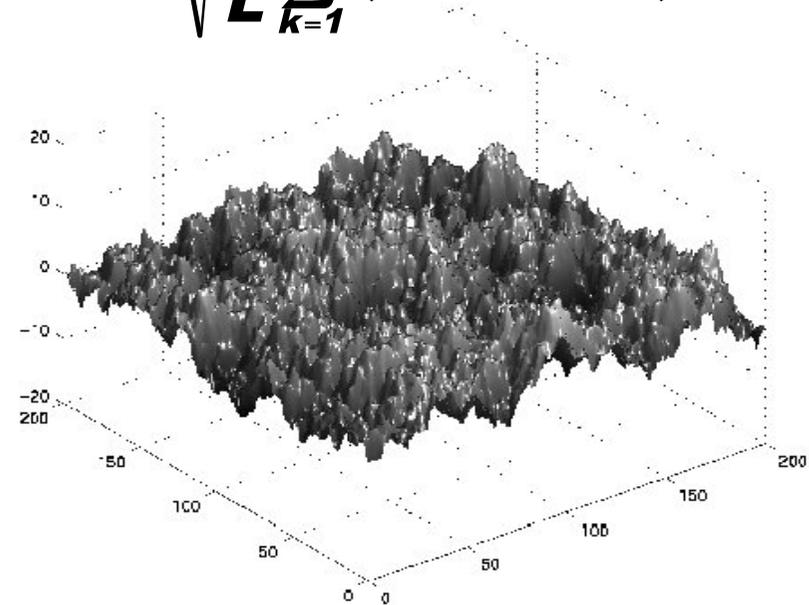
This *is* the **PDES** model

Non-equilibrium surface growth



$$\langle w^2(t) \rangle \sim \begin{cases} t^{2\beta}, & \text{if } t \ll t_x \\ L^{2\alpha}, & \text{if } t \gg t_x \end{cases}$$

$$w(\mathbf{t}) = \sqrt{\frac{1}{L} \sum_{k=1}^L (\tau_k(\mathbf{t}) - \bar{\tau}(\mathbf{t}))^2}$$



Dynamic scaling:

$$\alpha = \beta z$$

β growth exponent

z dynamic exponent

α roughness exponent

Coarse graining for the stochastic time surface evolution

Korniss, Toroczkai, Novotny, Rikvold, PRL '00

$$\tau_i(t+1) = \tau_i(t) + \eta_i(t) \Theta[\tau_{i-1}(t) - \tau_i(t)] \Theta[\tau_{i+1}(t) - \tau_i(t)]$$

- $\Theta(\dots)$ is the Heaviside step-function
- $\eta_i(t)$ iid exponential random numbers

•
•
•

$$\frac{\partial h(x,t)}{\partial t} = \nu \frac{\partial^2 h(x,t)}{\partial x^2} + \lambda \left[\frac{\partial h(x,t)}{\partial x} \right]^2 + D_{kpz} \eta(x,t)$$

Kardar-Parisi-Zhang
equation

$$P[\tau(x)] \propto \exp \left[-\frac{1}{2D} \int dx \left(\frac{\partial \tau}{\partial x} \right)^2 \right]$$

Steady state ($d=1$):
Edwards-Wilkinson
Hamiltonian

❖ Random-walk profile: short-range correlated local slopes

“Simulating the simulations”

❖ Universality/roughness (d=1)

$$\langle w^2(t) \rangle_L \sim \begin{cases} t^{2\beta}, & \text{if } t \ll t_x \\ L^{2\alpha}, & \text{if } t \gg t_x \end{cases}, \quad t_x \sim L^z, \quad z = \alpha / \beta$$

Foltin et al., '94

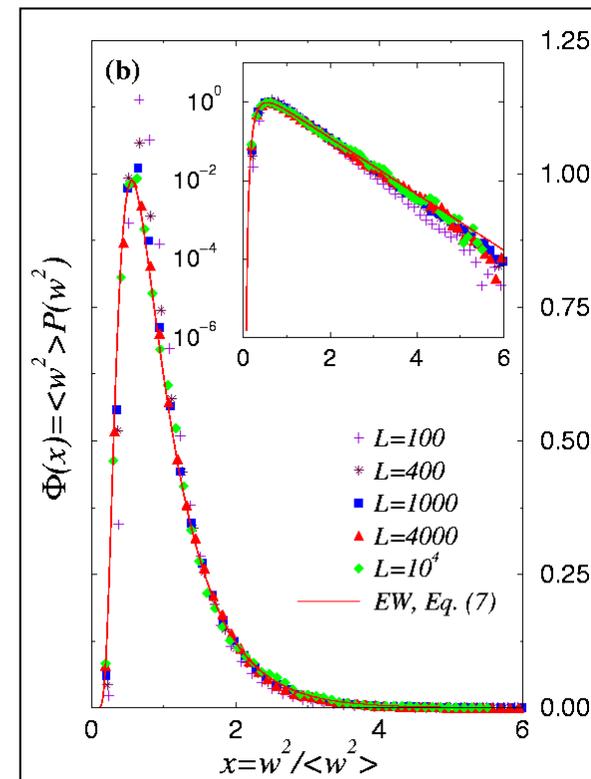
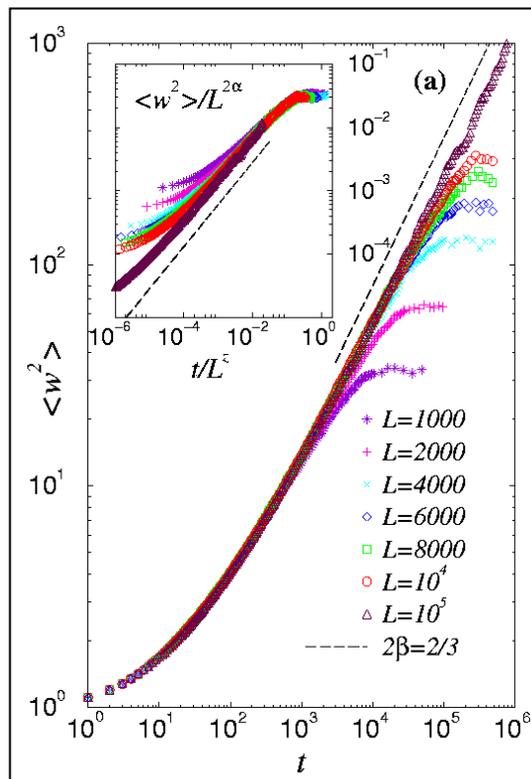
$$\beta \approx 0.33, \quad \alpha \approx 0.5$$

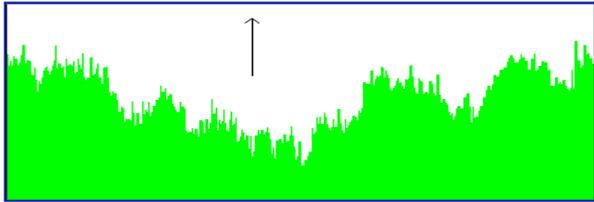
$$P(w^2) = \langle w^2 \rangle^{-1} \Phi(w^2 / \langle w^2 \rangle)$$

exact KPZ:

$$\beta = 1/3$$

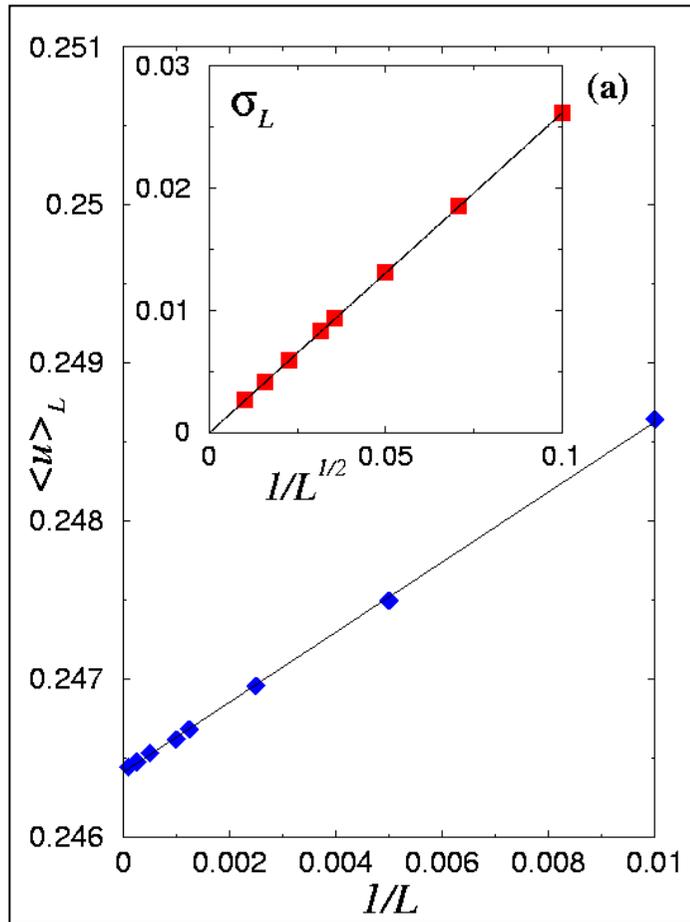
$$\alpha = 1/2$$





❖ Utilization/efficiency

Finite-size effects for the density of local minima/average growth rate (steady state):

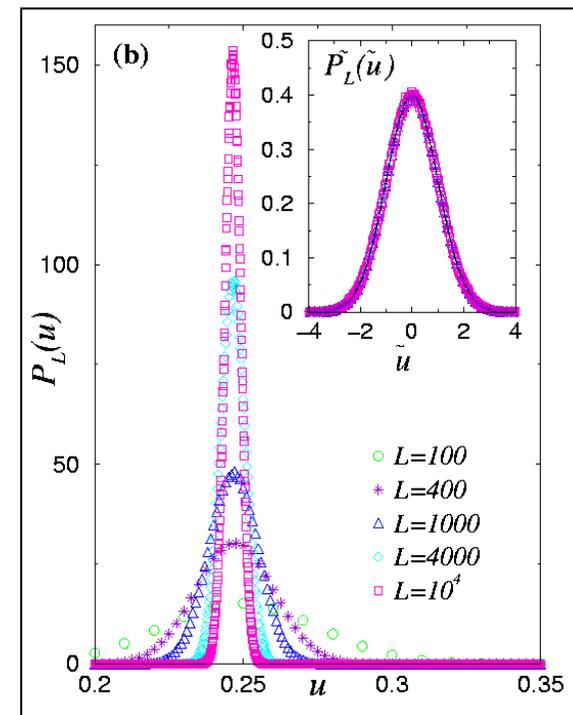


$d=1$

$$\langle u \rangle_\infty \approx 0.2464$$

$$\langle u \rangle = \langle u \rangle_\infty + \text{const}/N_{\text{PE}}$$

$$\sigma_L = \sqrt{\langle u^2 \rangle_L - \langle u \rangle_L^2} \sim 1/L^{1/2}$$



Implications for scalability

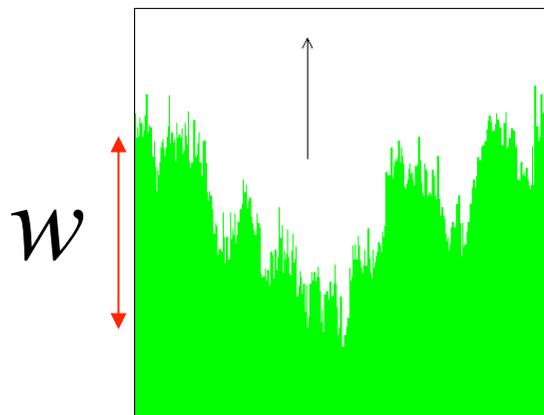
Virtual Time Horizon belongs to KPZ universality class

GREAT News ----- **Bad News**

❖ Simulation phase: **scalable** $\langle u \rangle = \langle u \rangle_\infty + \text{const}/N_{\text{PE}}$

$\langle u \rangle_\infty$ asymptotic average rate of progress of the simulation (utilization) is **non-zero**

❖ Measurement (data management) phase: **not scalable**



$$\langle w^2(t) \rangle_L \sim \begin{cases} t^{2\beta}, & \text{if } t \ll t_x \\ L^{2\alpha}, & \text{if } t \gg t_x \end{cases}$$

Rough Times!

Improve efficiency

Mixing

KPZ + RD

$$\frac{\partial h(x, t)}{\partial t} = \nu \frac{\partial^2 h(x, t)}{\partial x^2} + \lambda \left[\frac{\partial h(x, t)}{\partial x} \right]^2 + D_{kpz} \eta(x, t)$$

+

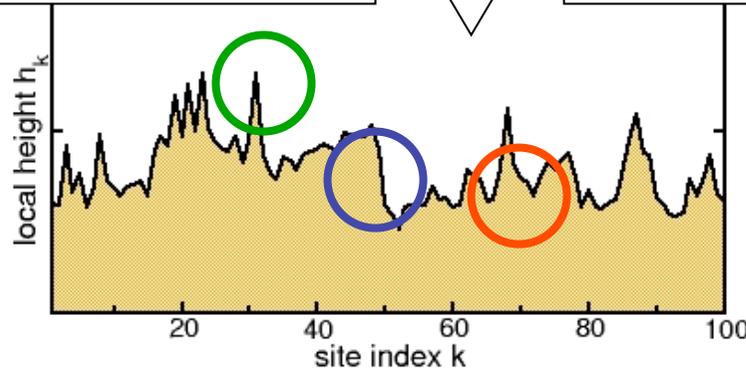
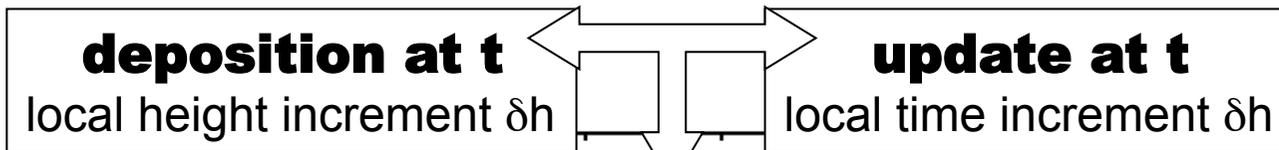
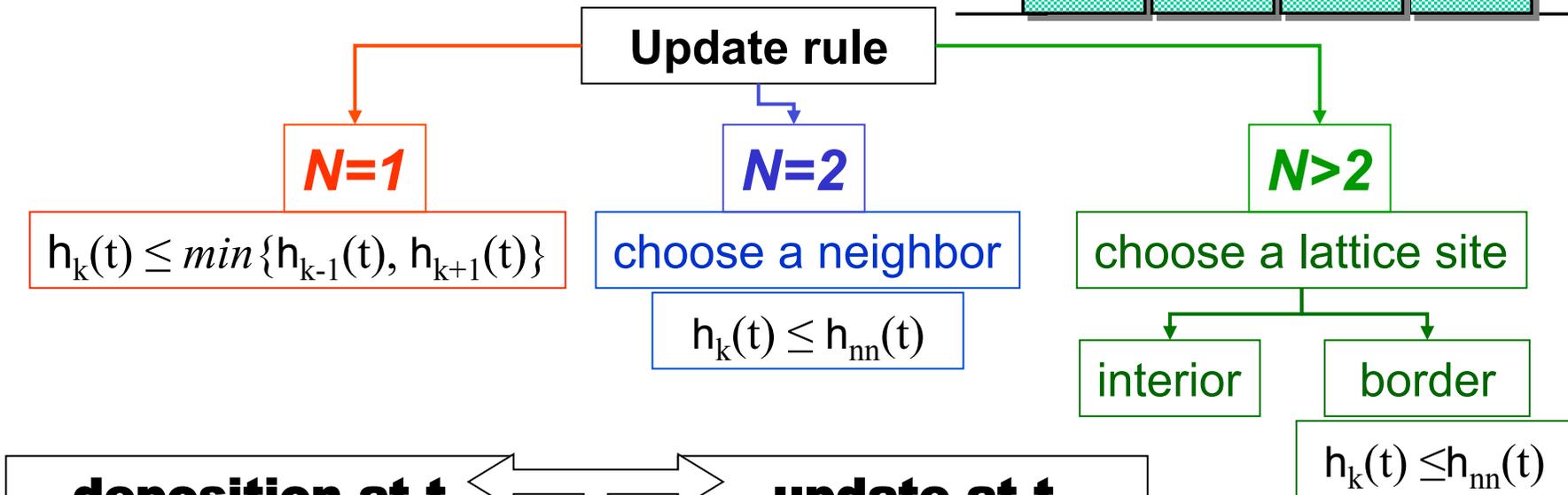
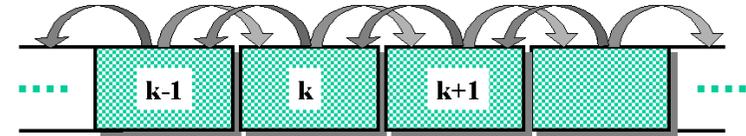
$$\frac{\partial h(x, t)}{\partial t} = D_{rd} \eta(x, t)$$

Simulation model for conservative PDES

Time-step t : index of the simultaneous update attempt

Updates at t : independent Poisson-random processes

If update at t : $h_k(t+1) = h_k(t) + \eta_{lk}(t)$

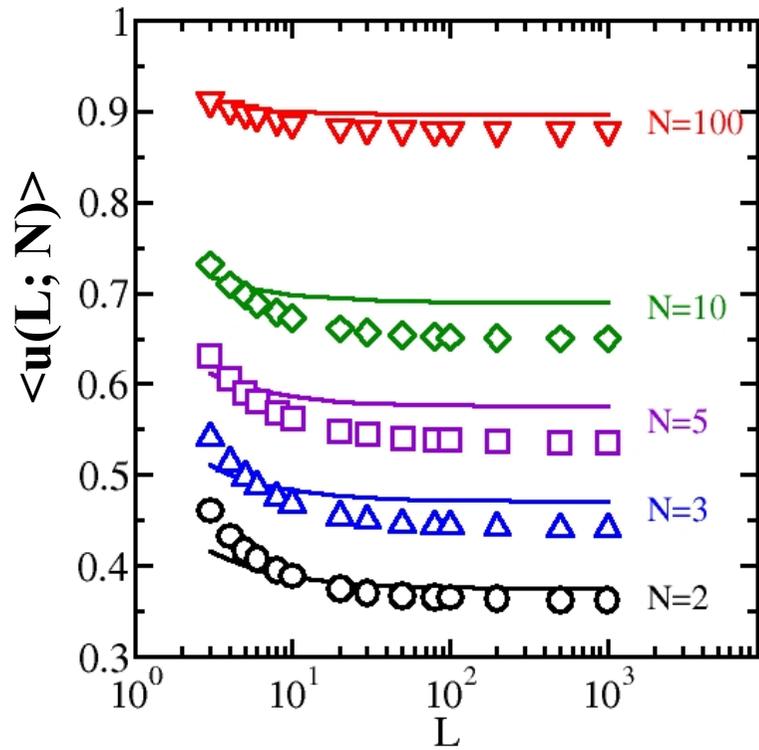


Virtual Time Horizon (VTH)

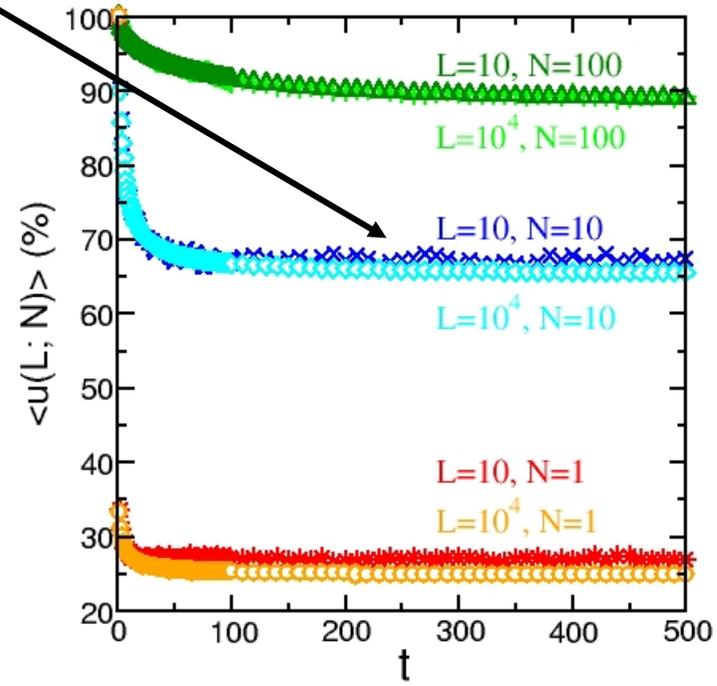
Properties of the algorithm are encoded in the VTH

Diagnostics: utilization of the parallel processing environment

Steady-state simulations



$$v(t) = \langle u(t) \rangle \mu_P$$

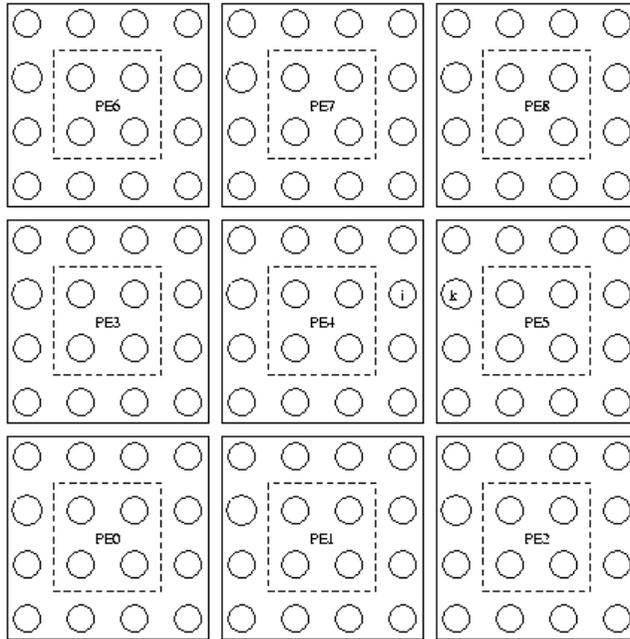


$$\begin{aligned} \langle u(L > 2; N > 1) \rangle &= \\ &= \left(1 - \frac{1}{\sqrt{2N}}\right) \left(1 - \frac{1}{2\sqrt{2N}} \frac{L-1}{L}\right) \end{aligned}$$

PRB 69, 075407 (2004)

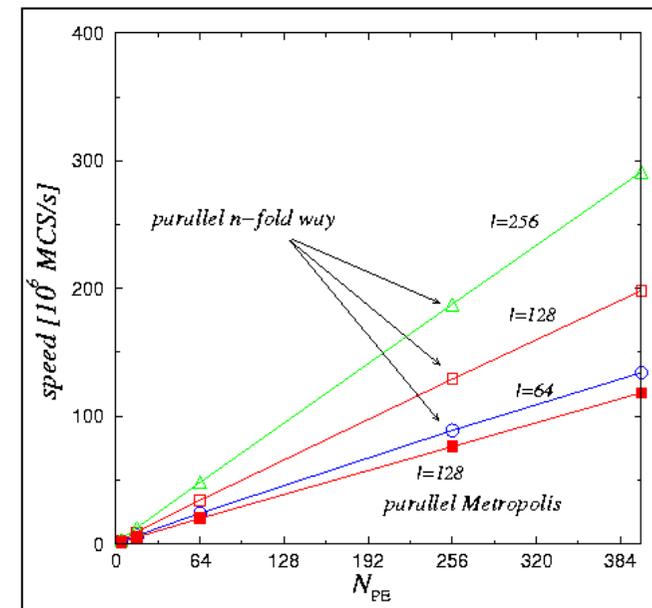
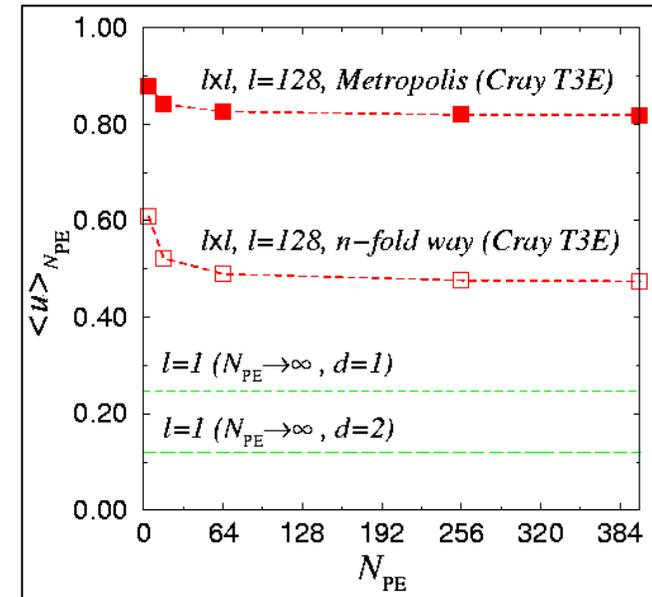
Actual implementation

Dynamics of a thin magnetic film



1. Local time incremented
2. Randomly chosen site
3. If chosen site is on the boundary, PE must wait until $\tau \leq \min\{\tau_{nn}\}$

$l > 1 \longrightarrow$ Mixing RD+KPZ



Implications for scalability

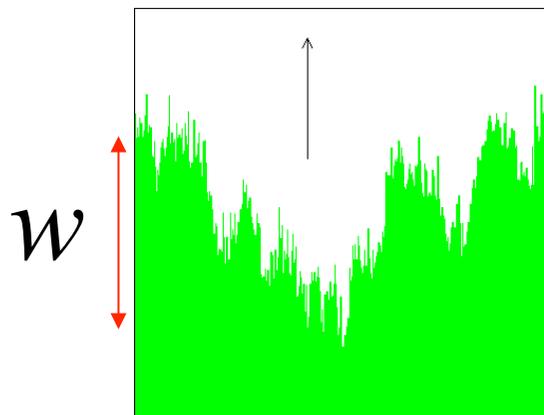
Virtual Time Horizon belongs to KPZ universality class

GREAT News ----- **Bad News**

❖ Simulation phase: **scalable** $\langle u \rangle = \langle u \rangle_\infty + \text{const}/N_{\text{PE}}$

$\langle u \rangle_\infty$ asymptotic average rate of progress of the simulation (utilization) is **non-zero**

❖ Measurement (data management) phase: **not scalable**



$$\langle w^2(t) \rangle_L \sim \begin{cases} t^{2\beta}, & \text{if } t \ll t_x \\ L^{2\alpha}, & \text{if } t \gg t_x \end{cases}$$

Rough Times!

How to make the measurement phase scalable as well?

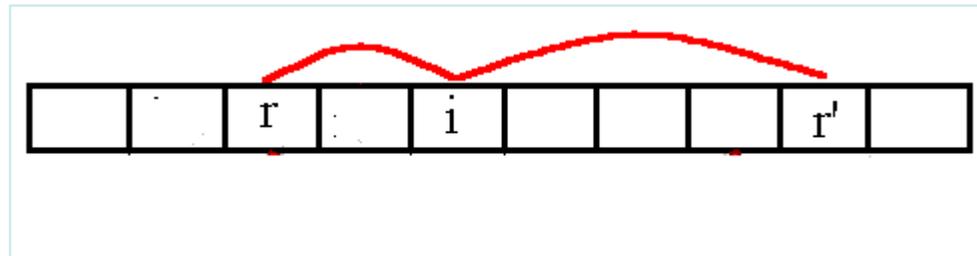
❖ Controlling the width

Greenberg et.al., '96

“Mean-field” like approximation to model the evolution of the time horizon (K-random interaction)

K=2: each PE **randomly** chooses two others, r and r'

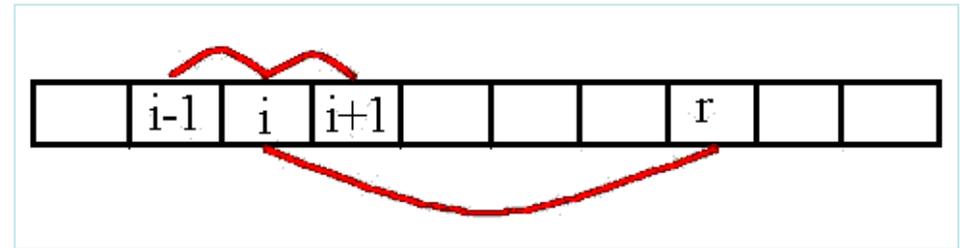
$$\tau_i \leq \min\{\tau_r, \tau_{r'}\}$$



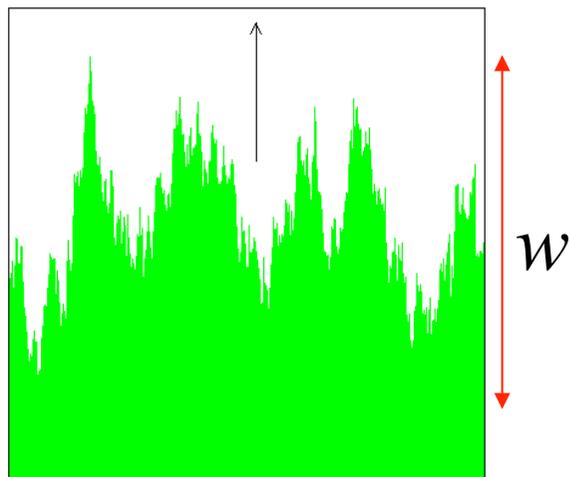
❖ Width is finite in this mean-field model when $L \rightarrow \infty$

❖ $\langle u \rangle_L \sim 1/(K+1)$ is nonzero

Annealed (or quenched) random connections

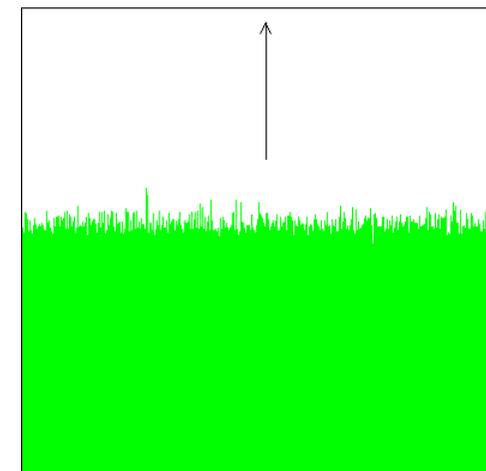


$$\tau_i \leq \min\{\tau_{nn}\}$$



KPZ surface: $w \sim L^\alpha$

$$\tau_i \leq \min\{\tau_{nn}, \tau_r\}$$

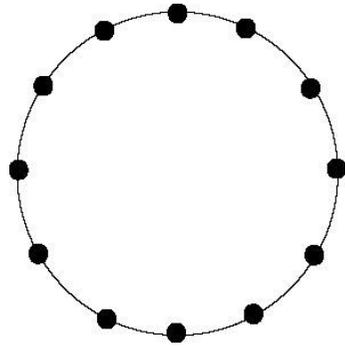


$w = \text{const.} + \mathbf{O}(L^{-1})$

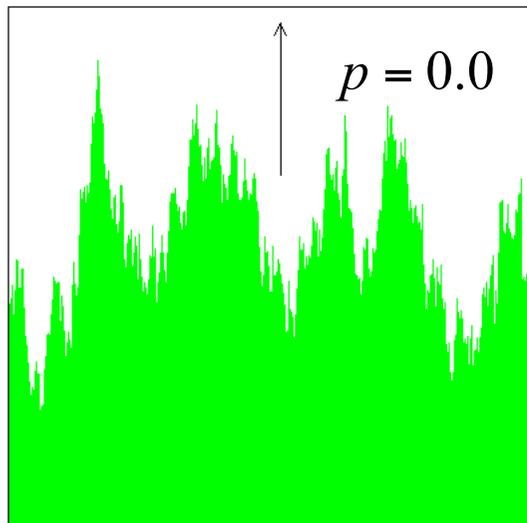
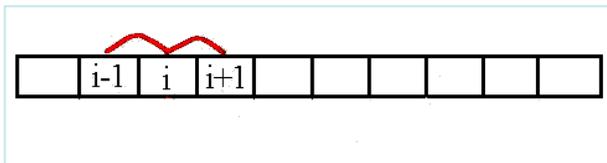
$$\partial_t \tau = \frac{\partial^2 \tau}{\partial x^2} - \lambda \left(\frac{\partial \tau}{\partial x} \right)^2 + \dots + \text{noise}$$

$$\partial_t \tau = -\gamma(\tau(x,t) - \bar{\tau}(t)) + \frac{\partial^2 \tau}{\partial x^2} + \dots + \text{noise}$$

Slopes are still short-range correlated: **non-zero** $\langle u \rangle$

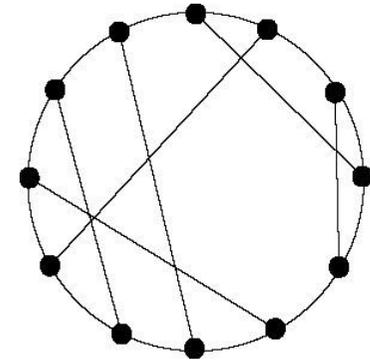


regular lattice (ring) topology
 (“ $p=0$ ”)

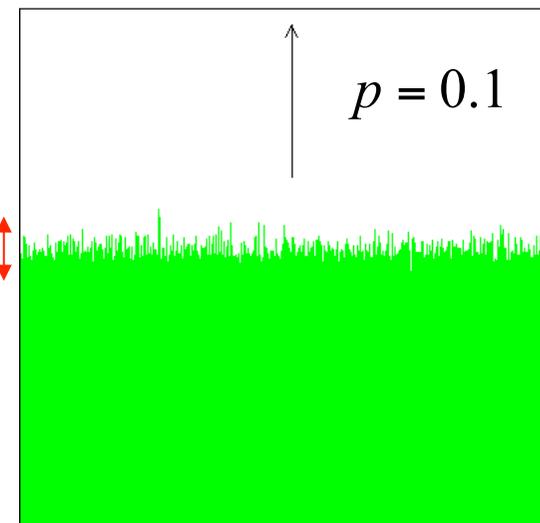
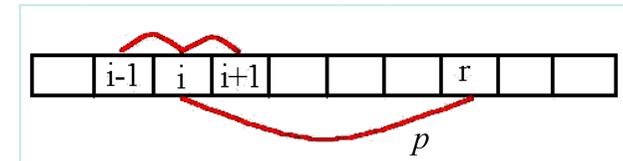


$$w \sim N^\alpha$$

$$N = 10^4$$



small-world-like connections:
 used with probability $p > 0$



$$N \rightarrow \infty$$

$$w \sim \text{const.}$$

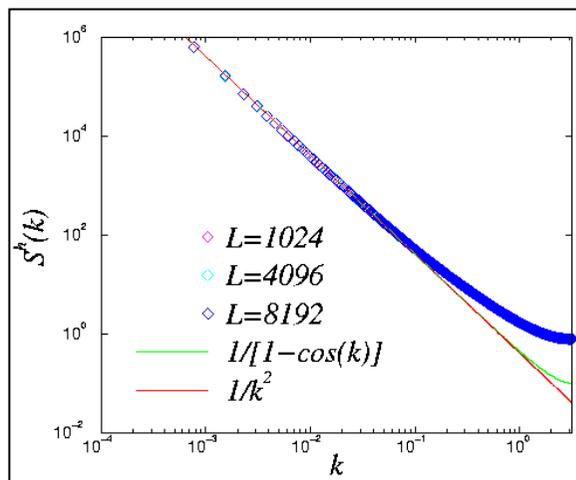
Steady-state “height” structure factors

$$S(k) \propto \langle \tau(k)\tau(-k) \rangle$$

only short-range connections (KPZ)

$$\partial_t \tau = \frac{\partial^2 \tau}{\partial x^2} - \lambda \left(\frac{\partial \tau}{\partial x} \right)^2 + \text{noise}$$

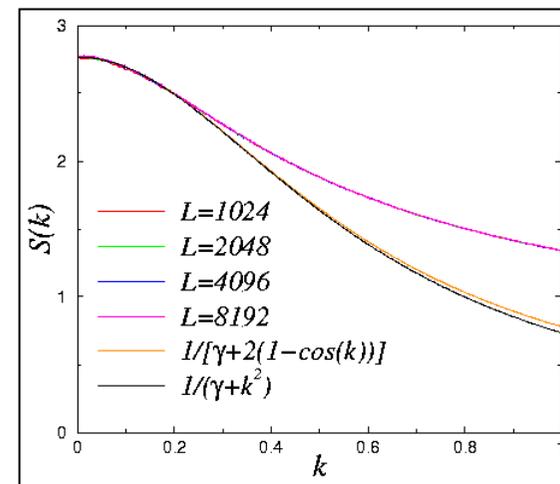
($d=1$) $S(k) \sim \frac{1}{k^2}$



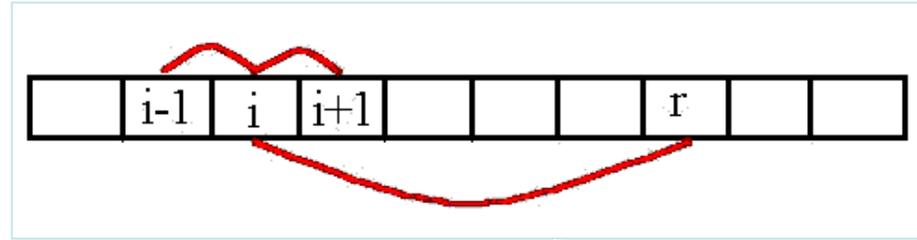
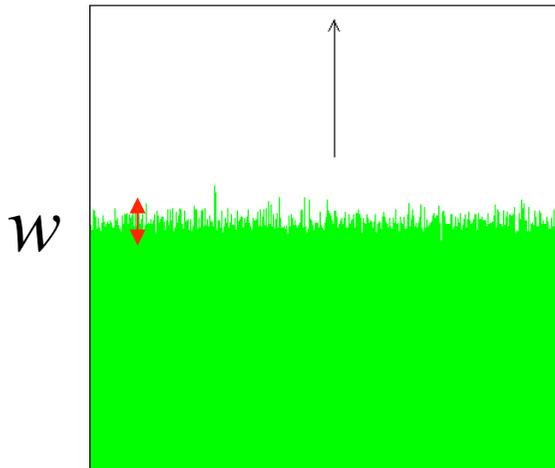
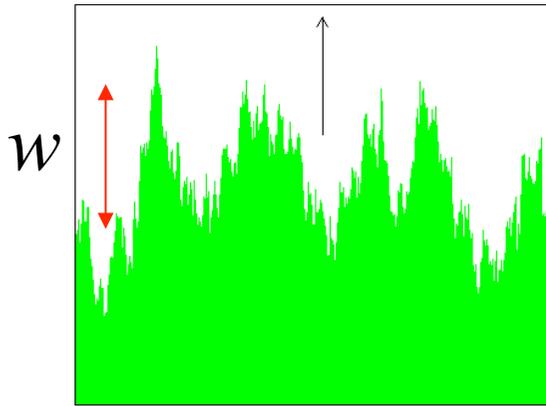
+ random connections (relaxation)

$$\partial_t \tau = -\gamma(\tau(x,t) - \bar{\tau}(t)) + \frac{\partial^2 \tau}{\partial x^2} + \text{noise}$$

$$S(k) \sim \frac{1}{\gamma + k^2}$$



Quenched random (Small World) connections



United States Patent
Novotny, et al.

6,996,504
February 7, 2006

Fully scalable computer architecture

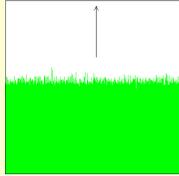
A scalable computer architecture capable of performing fully scalable simulations includes a plurality of processing elements (PEs) and a plurality of interconnections between the PEs. In this regard, the interconnections can interconnect each processing element to each neighboring processing element located adjacent the respective processing element, and further interconnect at least one processing element to at least one other processing element located remote from the respective at least one processing element. For example, the interconnections can interconnect the plurality of processing elements according to a fractal-type method or a quenched random method. Further, the plurality of interconnections can include at least one interconnection at each length scale of the plurality of processing elements.

Inventors: *Novotny*; Mark A. (Starkville, MS); Korniss; Gyorgy (Latham, NY)
Assignee: Mississippi State University (Mississippi State, MS)
Appl. No.: 990681
Filed: November 14, 2001

$$w = \text{const.} + \mathbf{O}(L^{-1}) \quad \tau_i \leq \min\{\tau_{nn}, \tau_r\}$$

Slopes are still short-range correlated: **non-zero** $\langle u \rangle$

PDES Summary and outlook

- Simple surface-growth model very useful
- The **tools and methods of non-equilibrium statistical physics** (coarse-graining, finite-size scaling, universality, etc.) can be applied to **scalability modeling and algorithm engineering**
- **Conservative** schemes can be made **perfectly scalable (ALL short-ranged PDES)**
 - Computational phase always scalable (KPZ universality)
 - Communication phase scalable with small-world network

Discussion and Provocations

- Neither software nor hardware nor algorithms alone will lead to (non-trivial) perfect scalability
- Without use of statistical mechanics, parallel computing will never be efficient/scalable
- Similar ideas apply to (non-trivial) cloud computing
- Similar ideas for sensor networks
- Similar ideas for databases and searches
- Similar ideas for fault-tolerant computing
- Similar ideas can be used to design new materials and devices with novel properties



1. **Synchronization in small-world-connected computer networks**, Guclu *et al*, PRE **73**, 066115 (2006).
2. **Universal scaling in mixing correlated growth with randomness**
Kolakowska *et al*, PRE **73**, 011603 (2006).
3. **Desynchronization and speedup in an asynchronous conservative parallel update protocol**, Kolakowska & Novotny, Ch. 6 in “Artificial Intelligence and Computer Science” (Nova Science 2005).
4. **Evolution of Time Horizons in Parallel and Grid Simulations**, L.N. Shchur and M.A. Novotny, PRE **70**, 026703 (2004).
5. **Roughening of the interfaces in (1+1) dimensional two-component surface growth with an admixture of random deposition**, Kolakowska *et al*, PRE **70**, 051602 (2004).
6. **Discrete-event analytic technique for surface growth problems**, Kolakowska and Novotny, PRB **69**, 075407 (2004).
7. **Suppressing roughness of virtual times in parallel discrete-event simulations**, Korniss, Novotny, Guclu, Toroczka, Rikvold, SCIENCE **299**, 677 (2003).
8. **Update statistics in conservative PDES**, Kolakowska *et al*, PRE **68**, 046705 (2003).
9. **Algorithmic scalability in globally constrained conservative PDES**
Kolakowska *et al*, PRE **67**, 046703 (2003).
10. **Algorithms for faster and larger dynamic Metropolis simulations**
Novotny *et al*, AIP Conference proceedings, 2003.
11. **Statistical Properties of the simulated time horizon in conservative PDES**
Korniss *et al*, ACM Proceedings, 2002.