

# Scientific Data Mining & Machine Learning

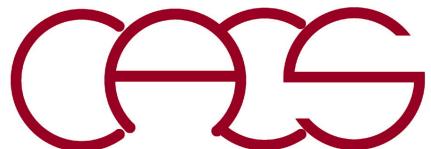
---

---

Aiichiro Nakano

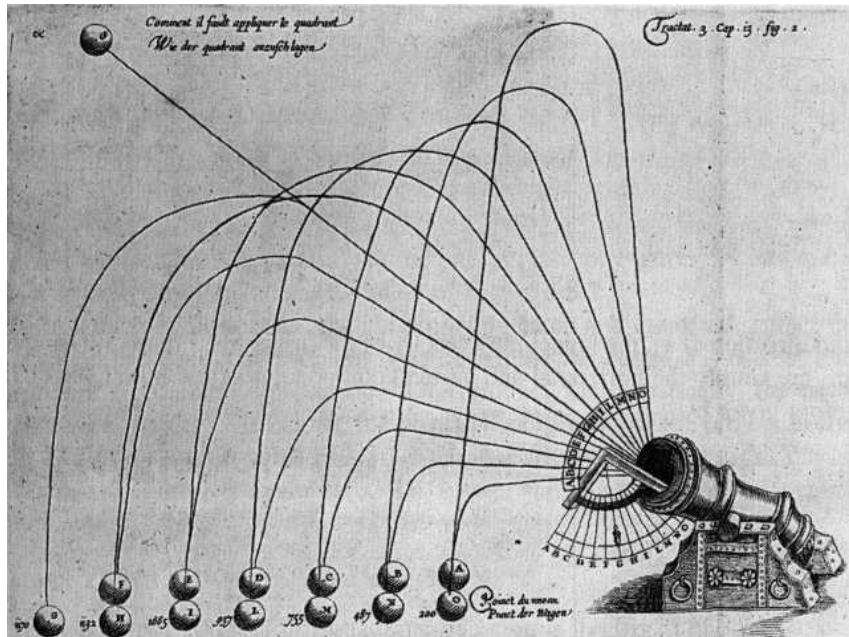
*Collaboratory for Advanced Computing & Simulations  
Department of Computer Science  
Department of Physics & Astronomy  
Department of Quantitative & Computational Biology  
University of Southern California*

Email: [anakano@usc.edu](mailto:anakano@usc.edu)



# Scientific Data Mining

- **Scientific data mining:** Automated detection of knowledge hidden in large & often noisy scientific (experimental, simulation, etc.) datasets
- **Knowledge:** Simplest (*i.e.*, minimal description length) explanation to replace exhaustive enumeration of the original data



Data



$$m \frac{d^2}{dt^2} \vec{r}(t) = \vec{F}$$

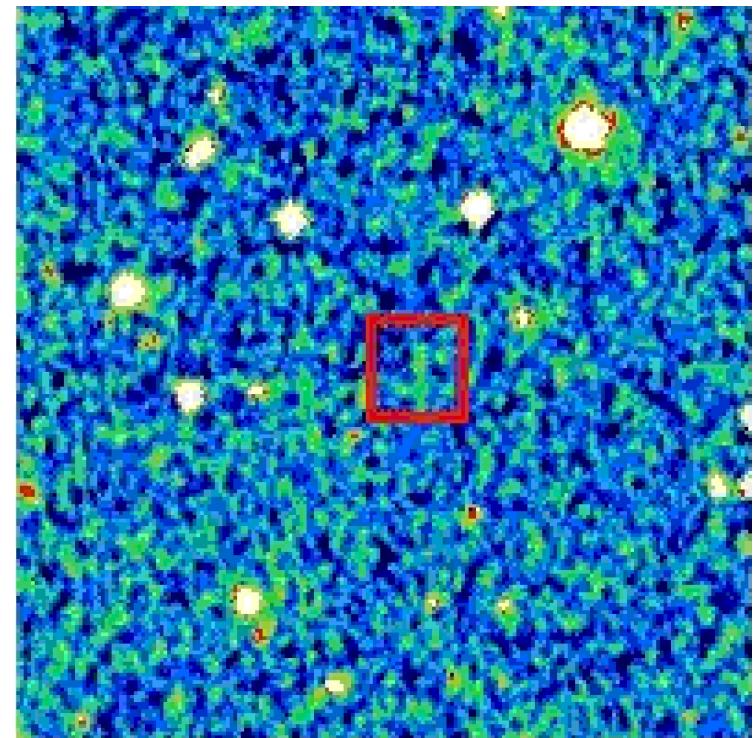
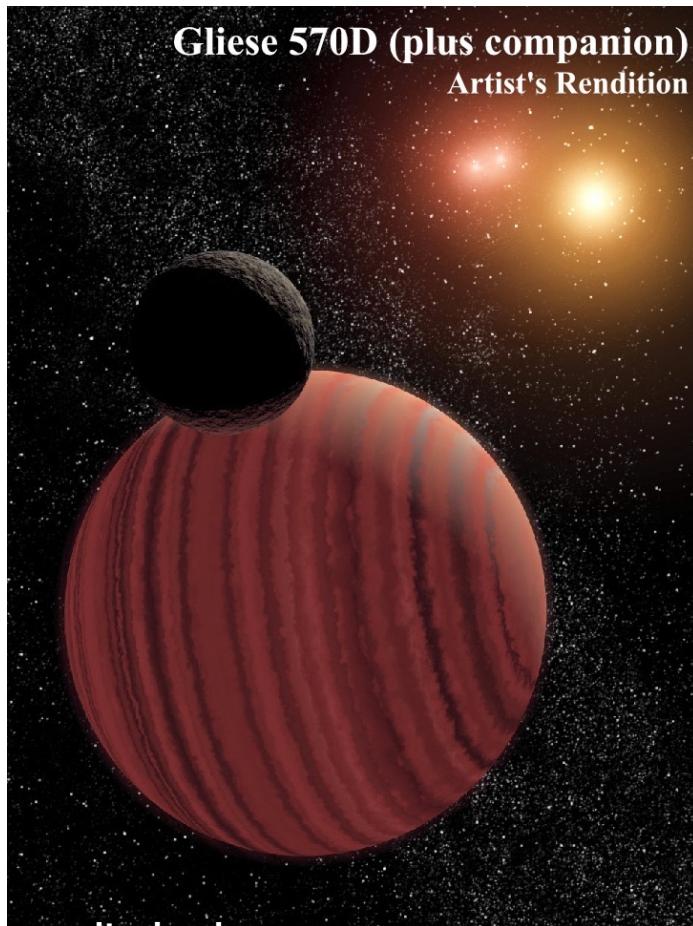
Knowledge

*cf.* [Occam's razor](#)  
[Kolmogorov complexity](#)

# Google Science in the Flat World

---

Parallel computing on globally distributed supercomputers & visualization platforms will revolutionize & democratize science & engineering (e.g., Google astronomy in the flat world)

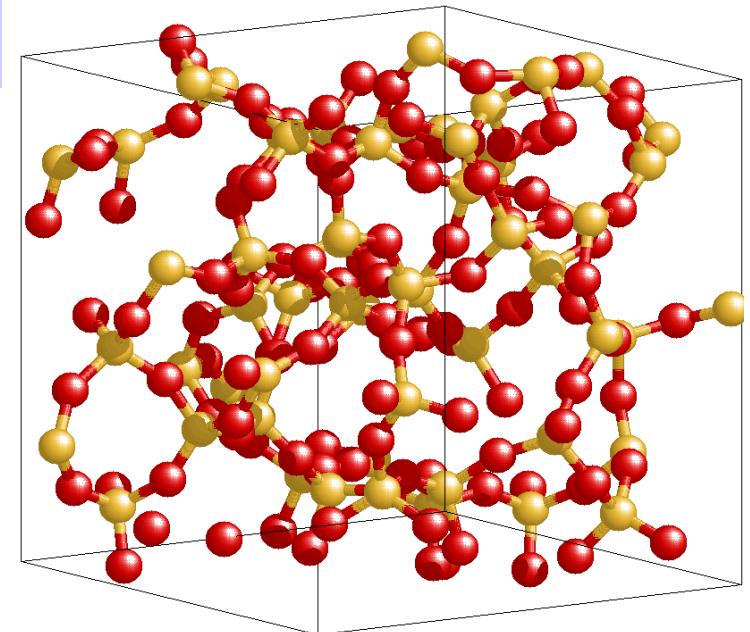


SDSS image of brown dwarf,  
2MASSI J0104075-005328

# Atomistic Data as a Graph

- Molecular dynamics data
  - Atomic data: species, positions, velocities, stresses, ...  
$$\{\lambda_i, \vec{r}_i, \vec{v}_i, \vec{\sigma}_i, \dots | i = 1, \dots, N\}$$
  - Atomic-pair data: bond order, pair distance, ...  
$$\{B_{ij}, \vec{r}_{ij}, \dots | i, j = 1, \dots, N; i \neq j\}$$
- Chemical bond network  $G = (V, E)$ 
  - Node degrees
  - Paths
  - Rings
  - Frequently occurring subgraphs

V: Set of atoms  
E: Set of bonds

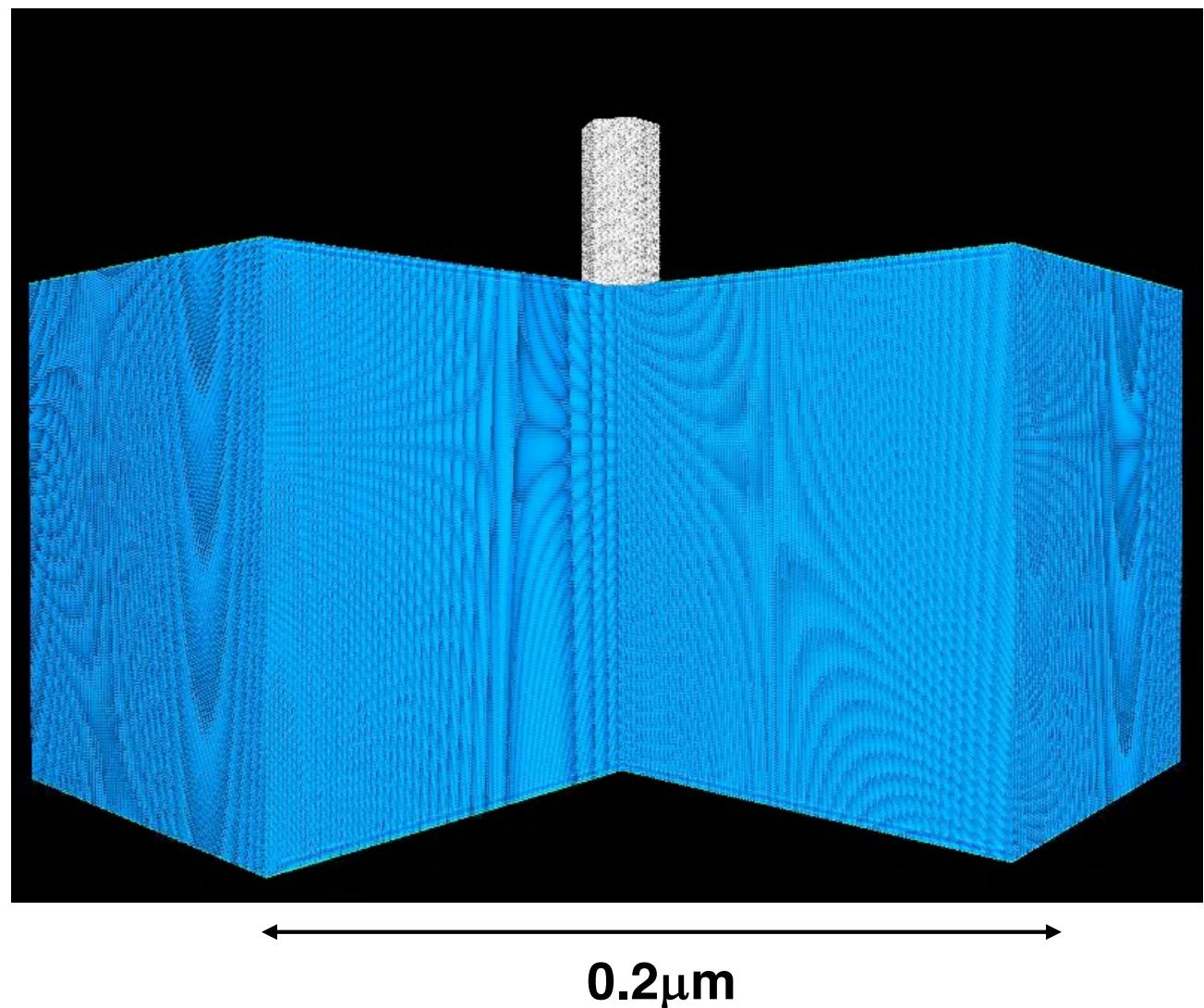


# Hypervelocity Impact on Ceramics

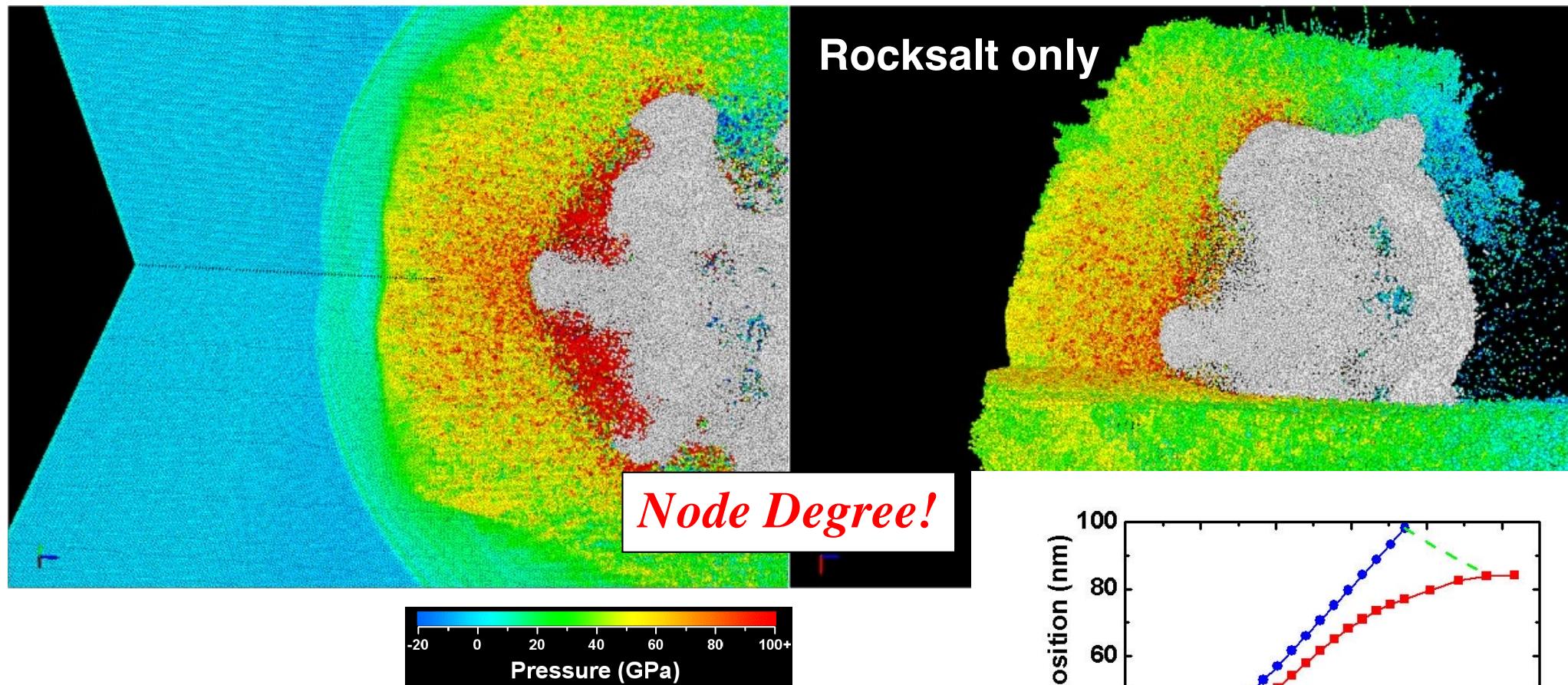
- 209M-atom MD of AlN
- 300M-atom MD of SiC
- 540M-atom MD of  $\text{Al}_2\text{O}_3$

↑ [0001]

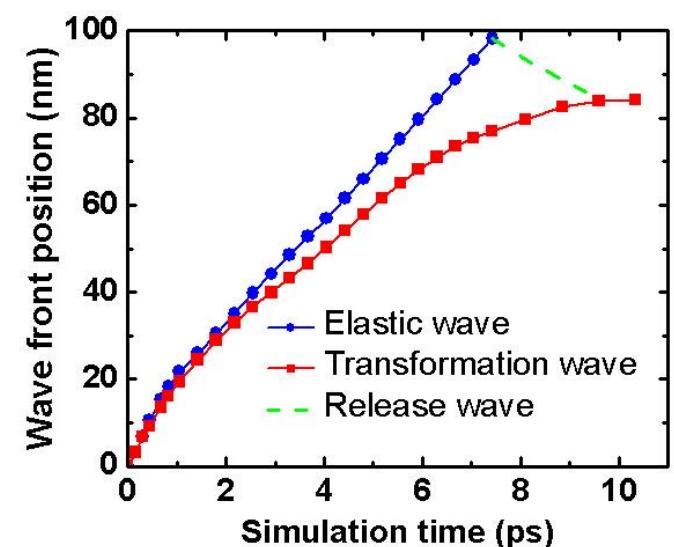
- $\text{Al}_2\text{O}_3$  plate
- 18 km/s impact



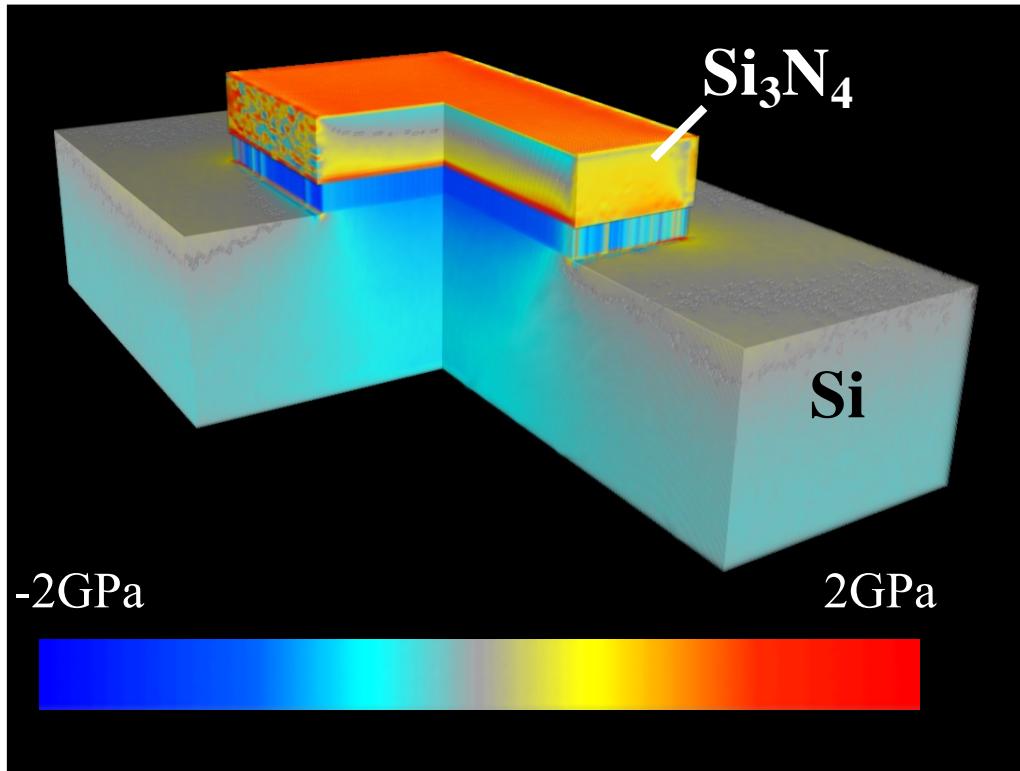
# Shock-Induced Structural Phase Transformation in AlN



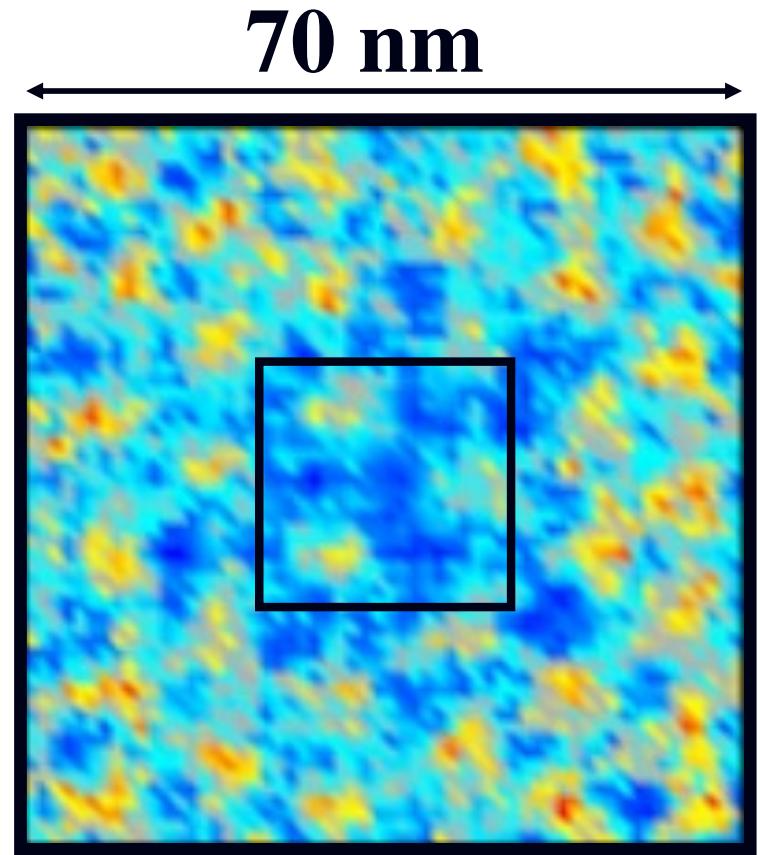
- Wurtzite (4-coordinated) to rocksalt (6-coordinated) phase transformation at 20 GPa



# Stress Domains in $\text{Si}_3\text{N}_4$ /Si Nanopixels



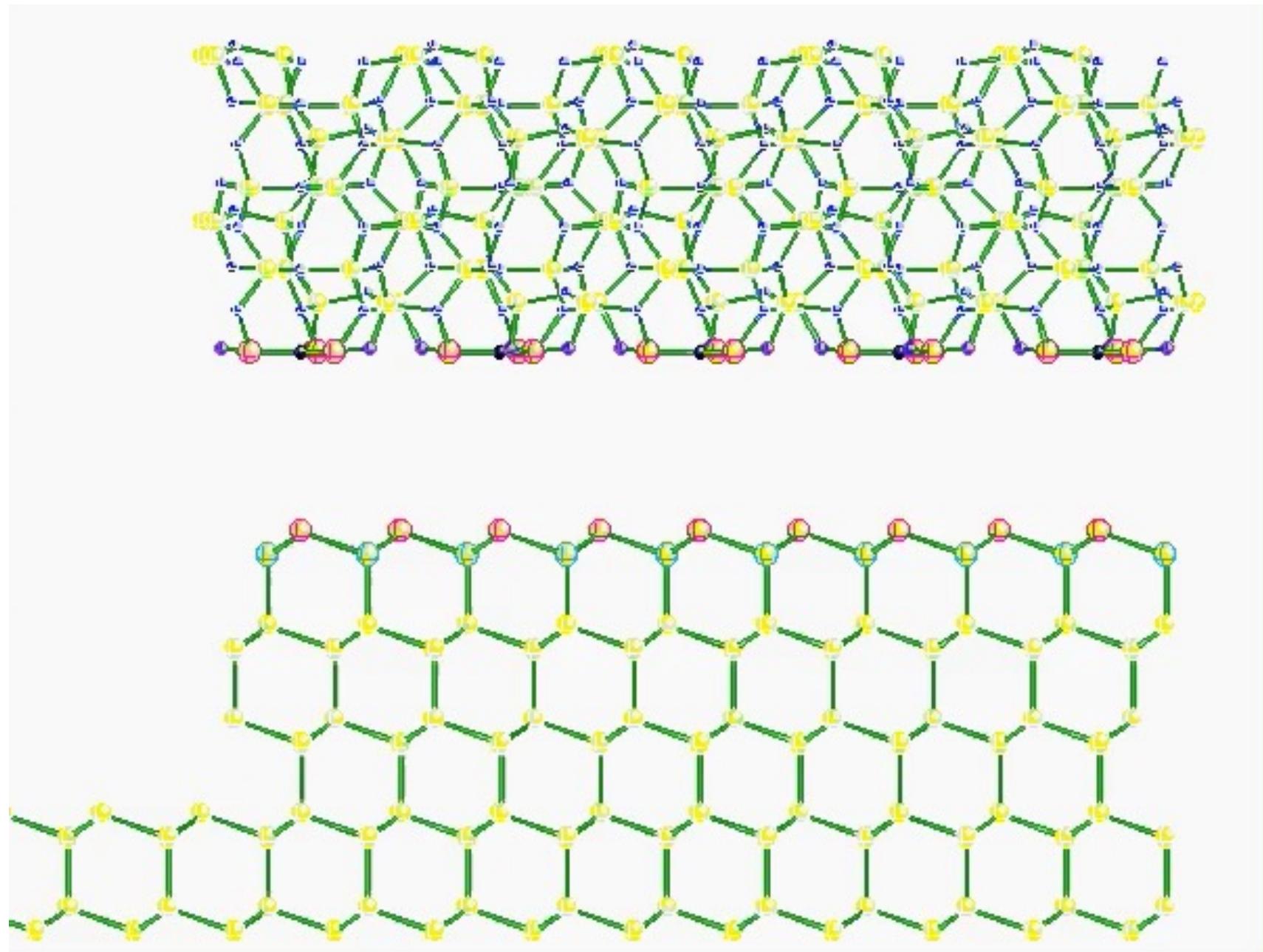
Stress well in Si with a  
crystalline  $\text{Si}_3\text{N}_4$  film  
due to lattice mismatch



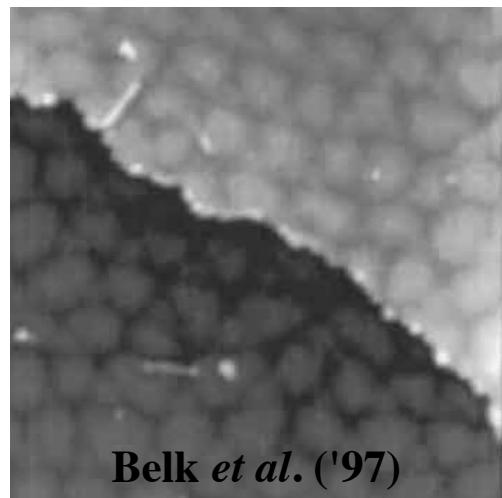
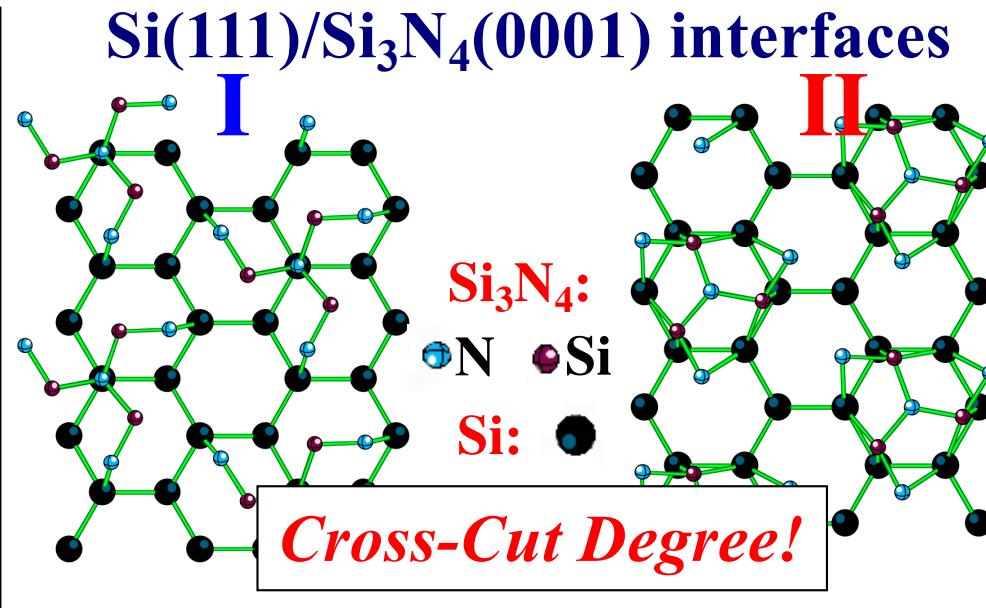
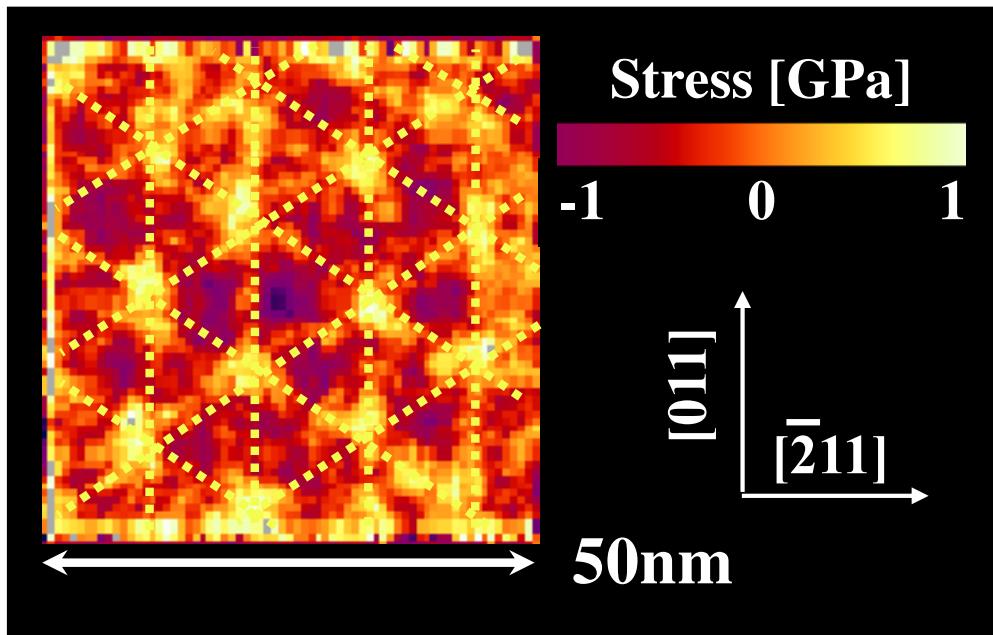
Stress domains in Si  
due to an amorphous  
 $\text{Si}_3\text{N}_4$  film

# **Si(111)/Si<sub>3</sub>N<sub>4</sub>(0001) Interface**

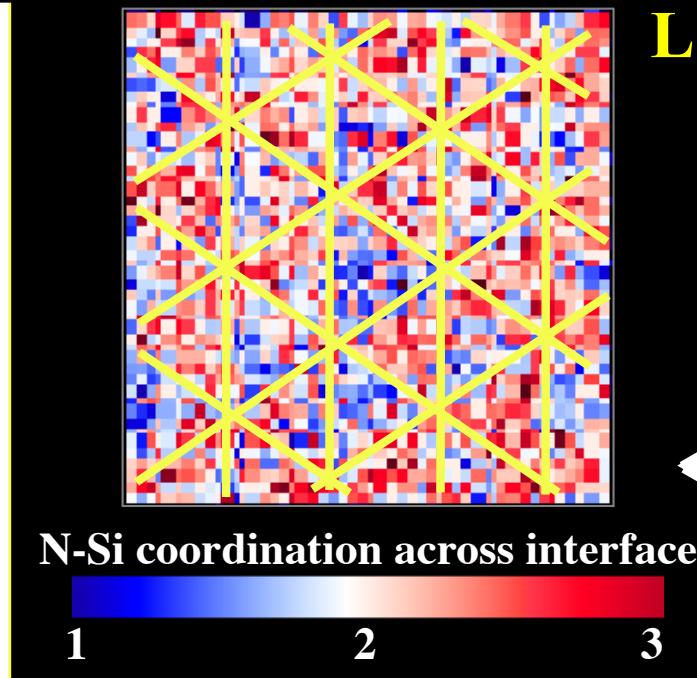
---



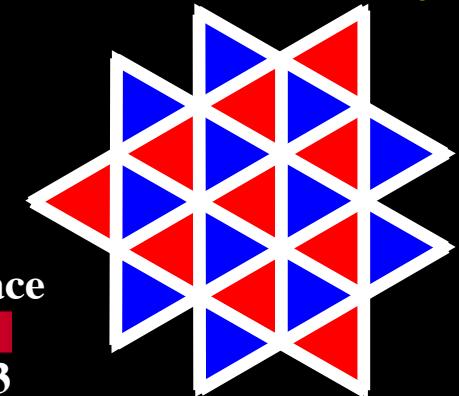
# Stress Domains in Si/Si<sub>3</sub>N<sub>4</sub> Nanopixel



Misfit dislocation network  
in InAs/GaAs(111)

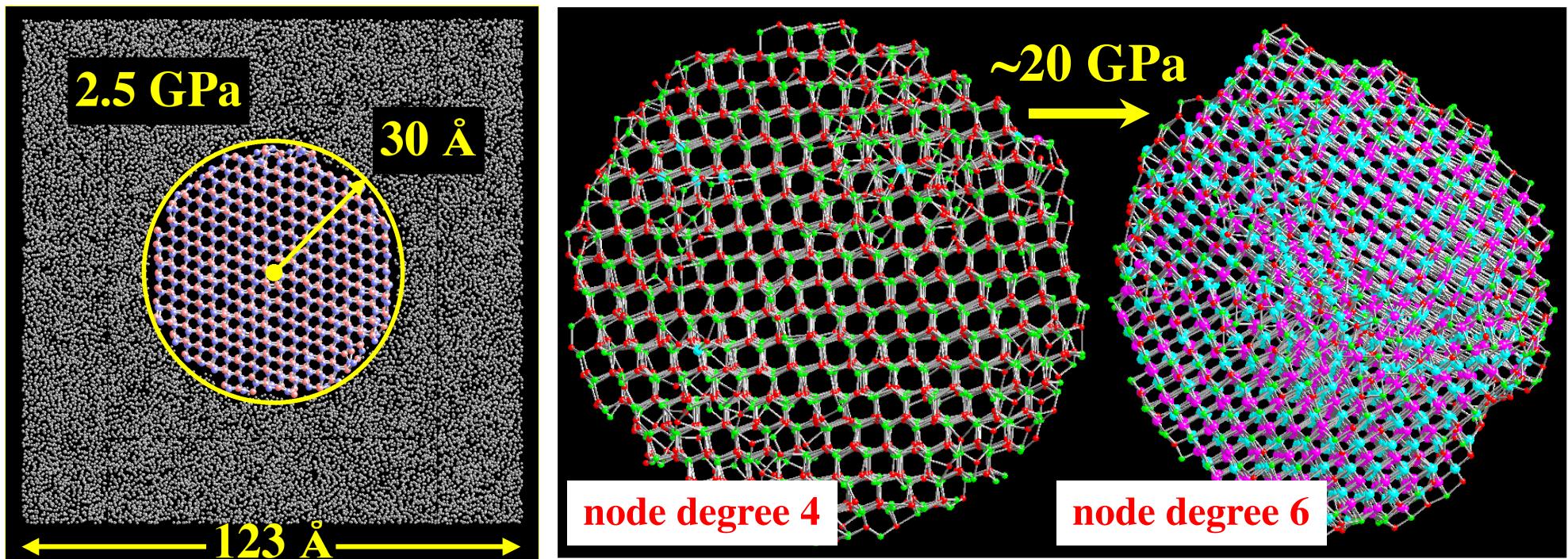


Lattice mismatch  
(1%) induced  
interfacial  
domain array



# High-Pressure Structural Transformation

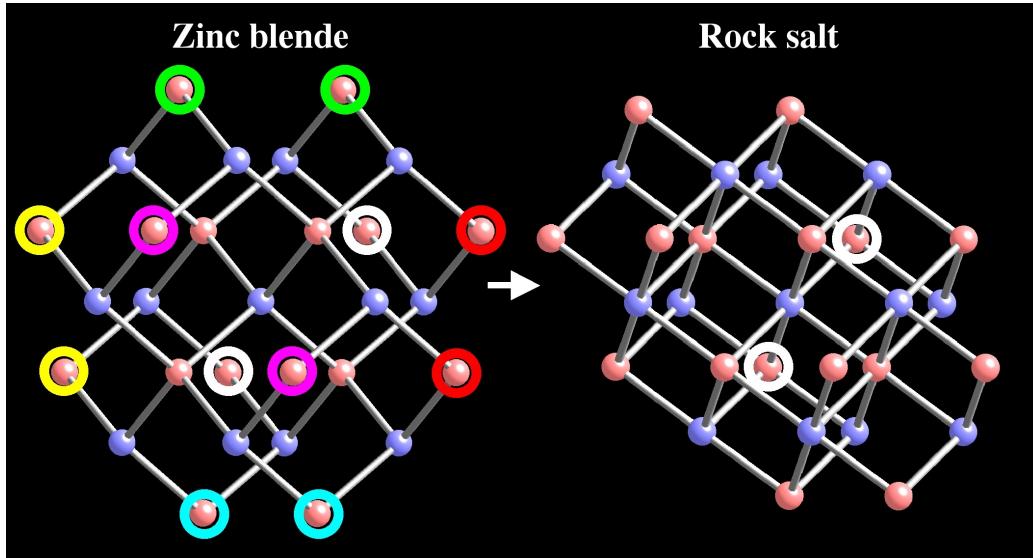
- Wurzite (node degree 4) to rocksalt (node degree 6) structural transformation of a GaAs nanoparticle under high pressure



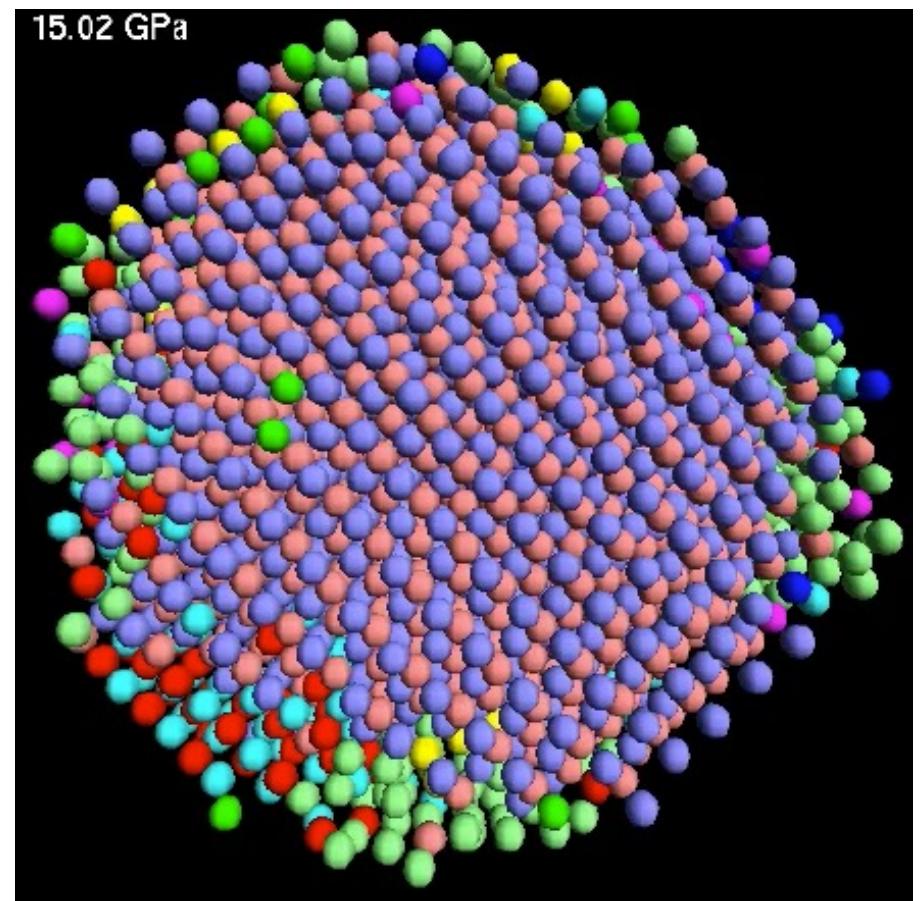
- Existence of multiple domains?

# Graph-Transition Tracking

- Finite set of graph transitions as a classifier



$$\begin{array}{c} G = (V, E) \\ \downarrow \\ G' = (V, E') \\ E \subset E' \end{array}$$



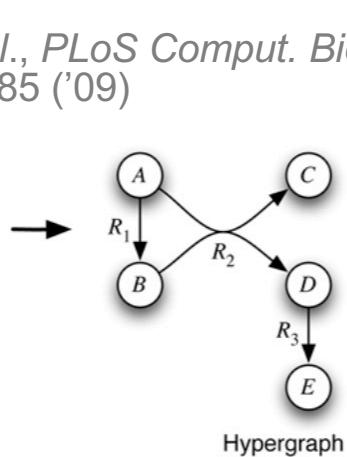
*Graph Transition!*

# Chemical Reaction Network

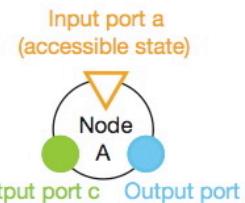
Klamt et al., PLoS Comput. Biol.  
5, e1000385 ('09)

Reaction networks

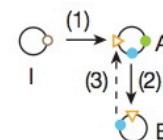
$$\begin{aligned} R_1 : A &\longrightarrow B \\ R_2 : A + B &\longrightarrow C + D \\ R_3 : D &\longrightarrow E \end{aligned}$$



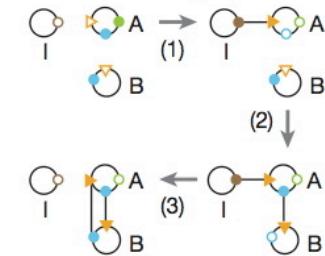
## c Nodal abstraction



## d Reaction graph

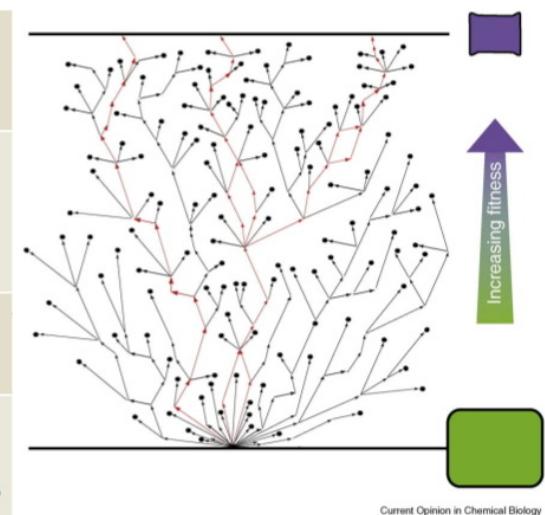
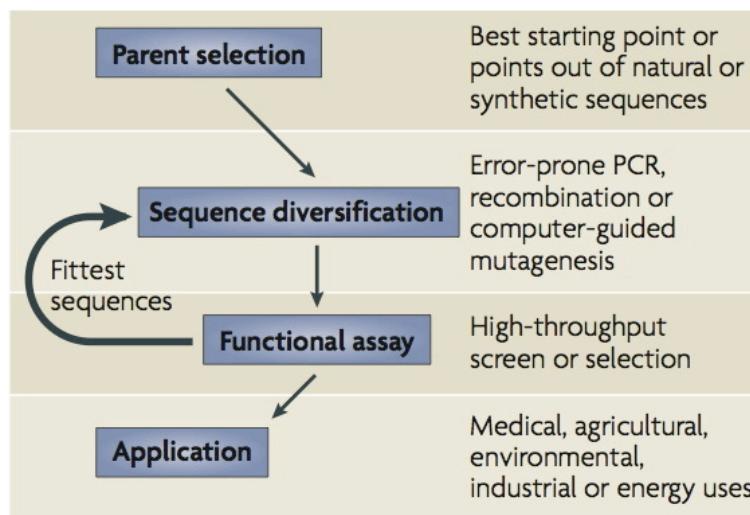
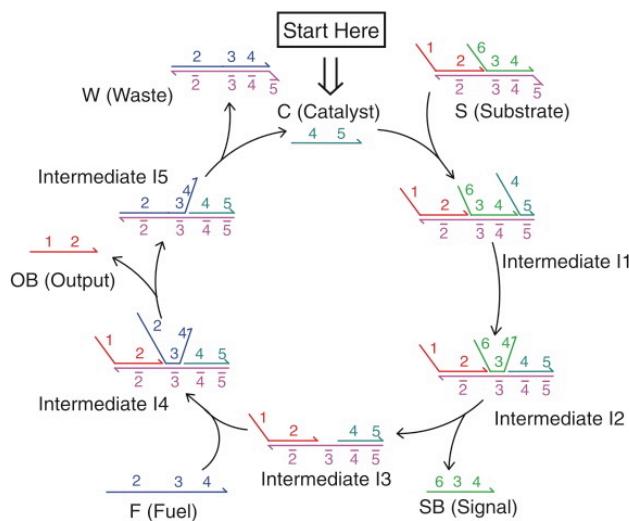


## e Execution of reaction graph



Yin et al., Nature 451, 318 ('08)

Zhang et al., Science 318, 1121 ('07)

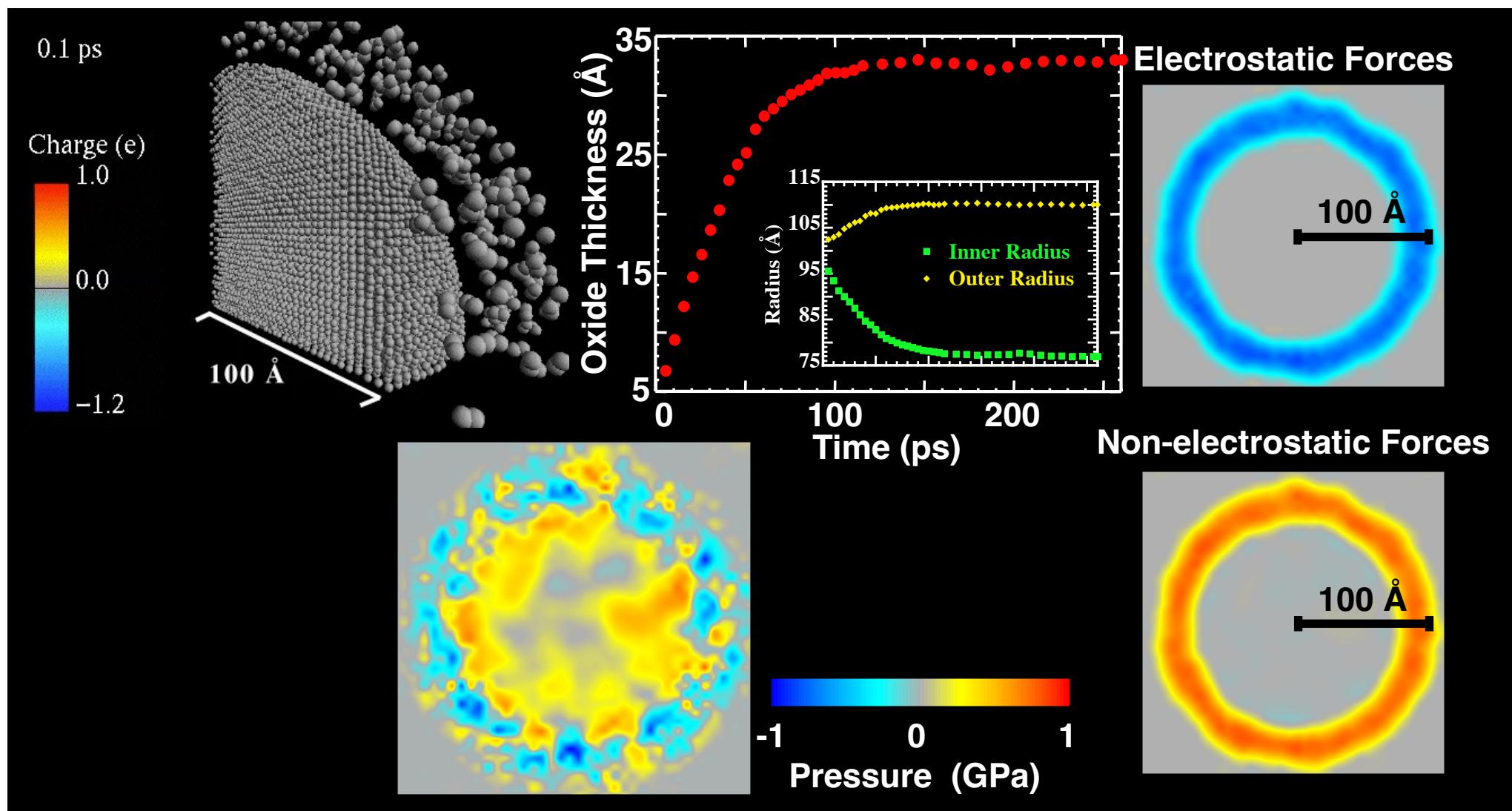


Arnold group, Nature Rev. MCB 10, 867('09); COCB 13, 3 ('09)

**Reaction graph = language for self-assembly & Directed & accelerated evolution catalytic cycle design**

Chen et al., Nature Nanotechnol. 8, 755 ('13)

# Oxidation of an Al Nanoparticle (n-Al)

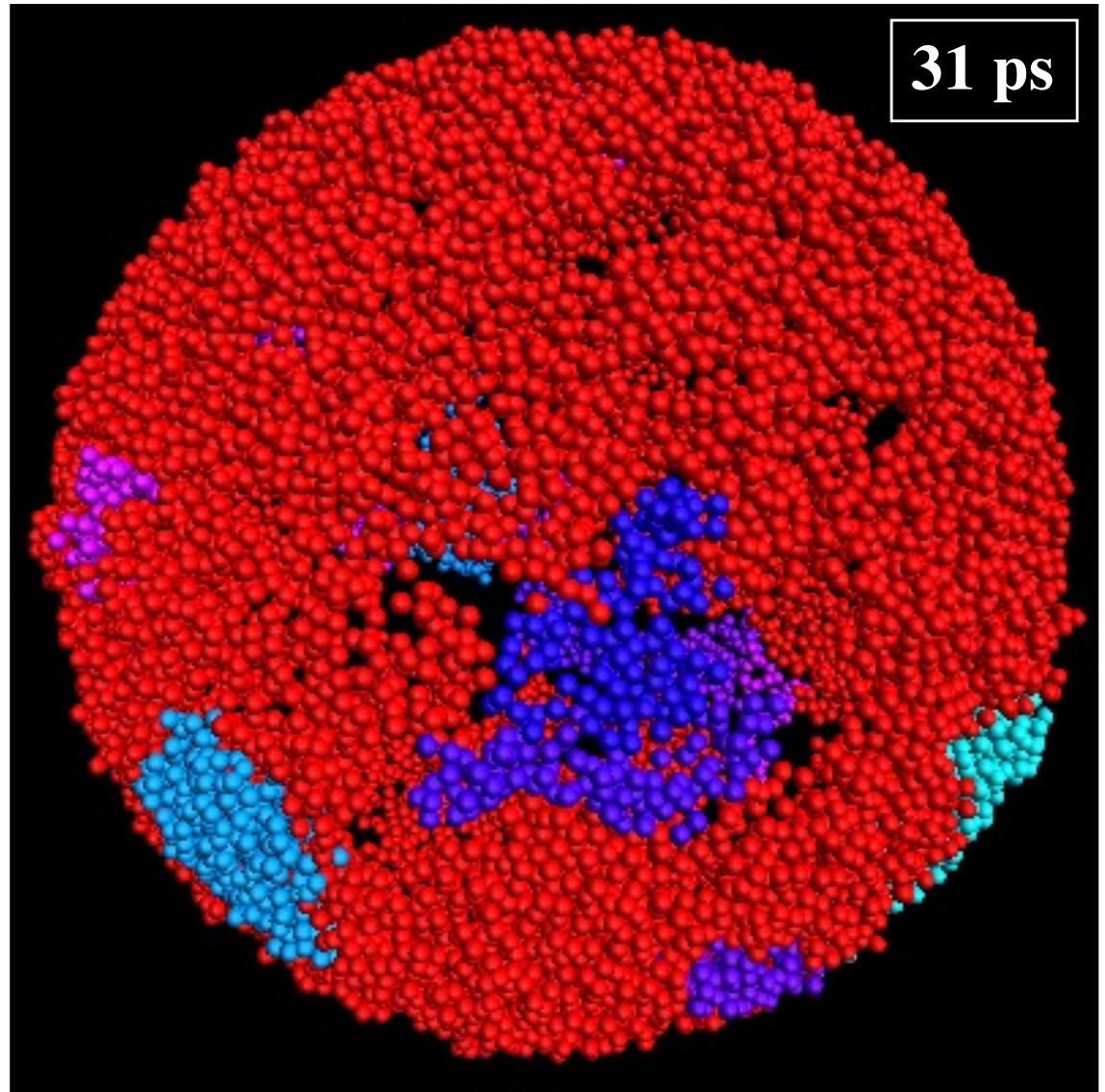


- Oxide thickness saturates at 40 Å after 0.5 ns, in agreement with experiments
- Oxide region/metal core is under negative/positive pressure
- Attractive Al-O Coulomb forces contribute large negative pressure in the oxide

# Oxidative Percolation

Clusters of  $\text{OAl}_4$  coalesce to form a neutral, percolating tetrahedral network that impedes further growth of the oxide

*Percorative  
Connected Components!*



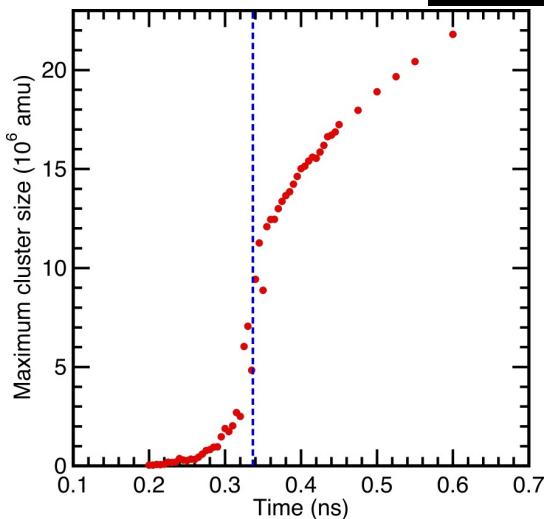
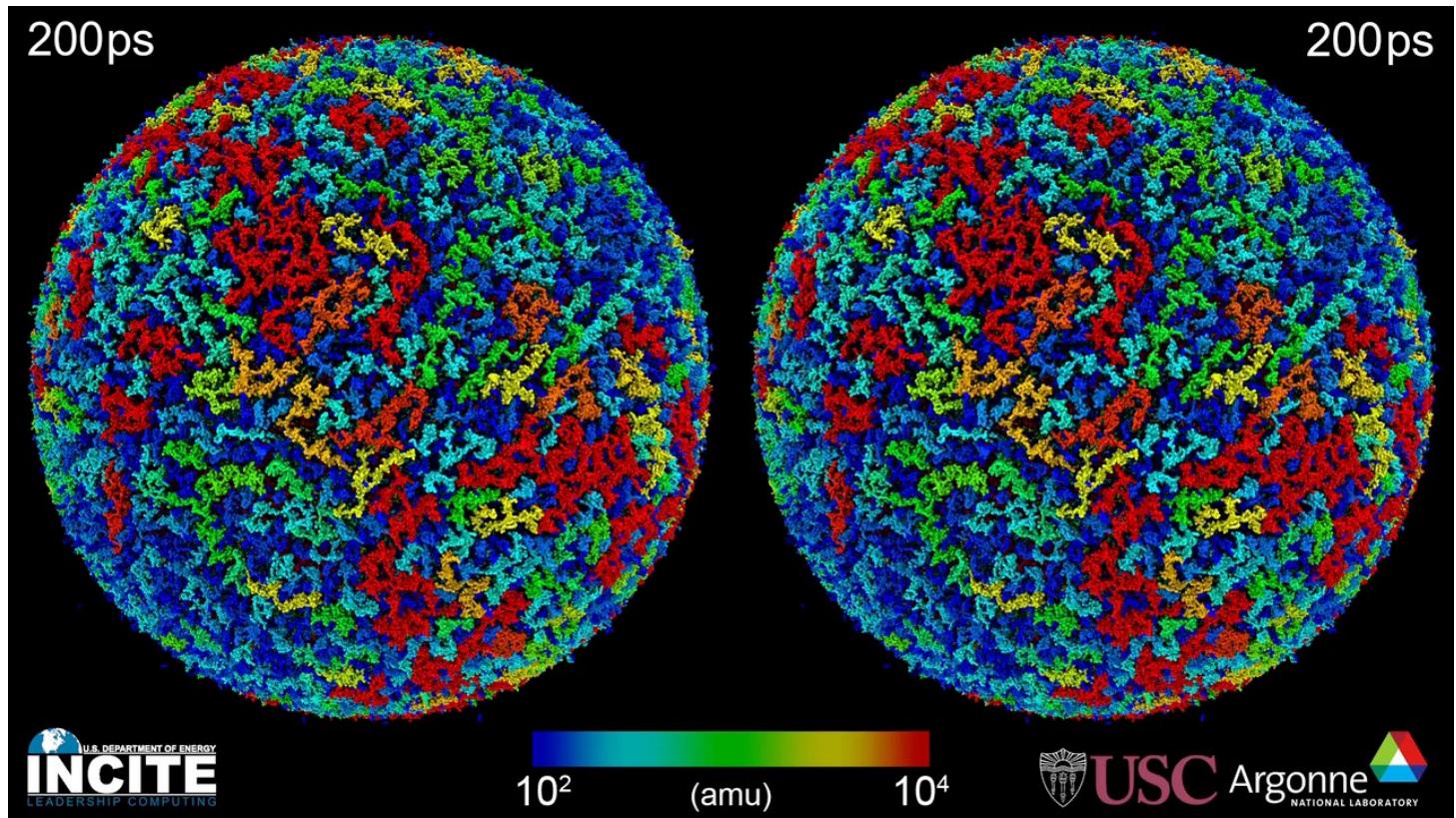
Size of Network

$10^2$        $10^3$        $10^4$

# Fractal Nanocarbon Product

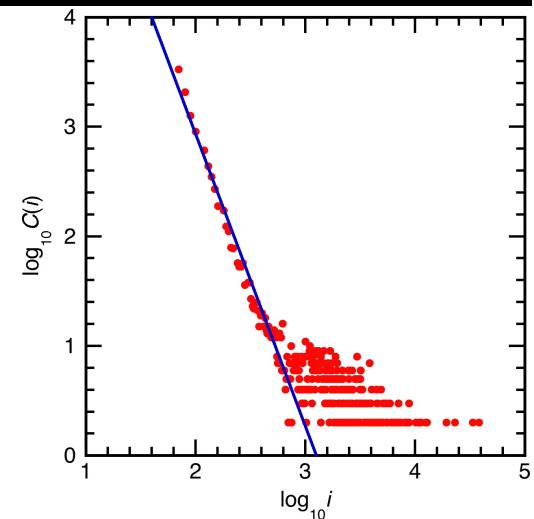
## Oxidation of SiC particle

- Percolation transition causes carbon clusters to exhibit power-law distribution of sizes:  $C(i) \sim i^{-\tau}$



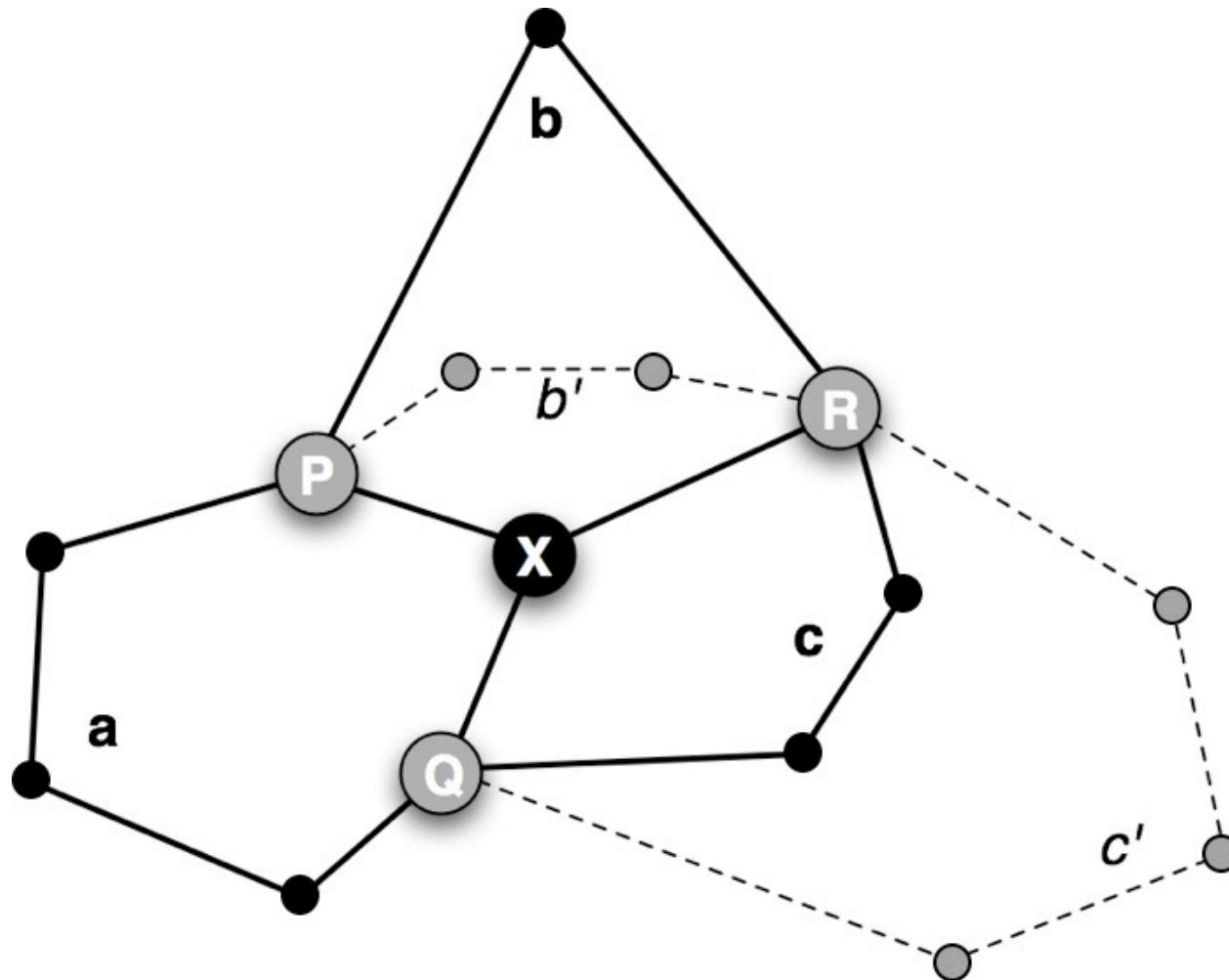
- Fractal nanocarbon product with large surface areas may find supercapacitor, battery-electrode & mechanical metamaterial applications:  $d_f = d/(\tau - 1) \sim 1.85$

K. Nomura *et al.*, *Sci. Rep.* **6**, 24109 ('16)  
J. Insley *et al.*, *IEEE/ACM SC16*



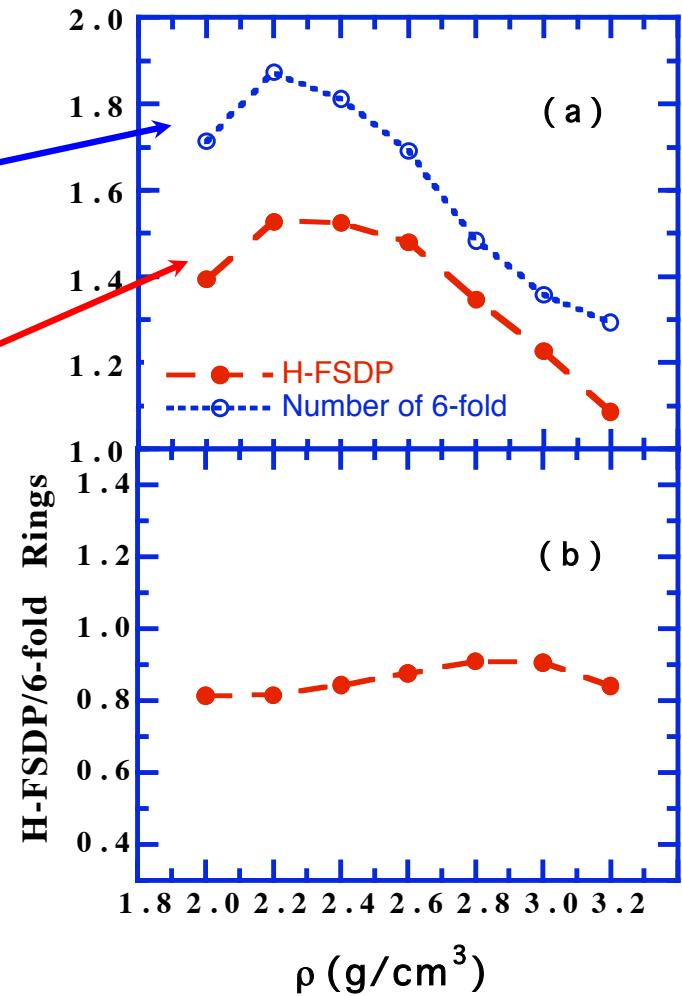
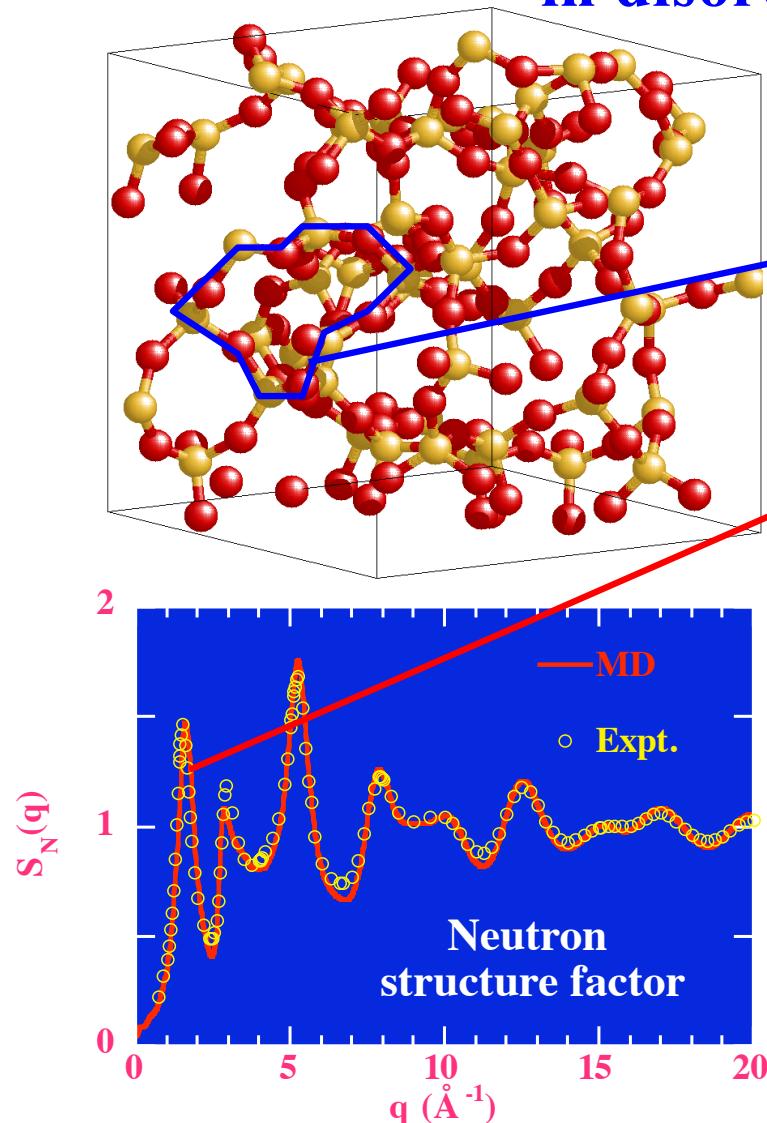
# Shortest-Path Rings

- **K-ring:** Given a vertex  $x$  & two of its neighbors  $w$  &  $y$ , a K-ring generated by the triplet  $w-x-y$  is any ring containing the edges  $[w-x]$ ,  $[x-y]$  and a shortest path  $w-y$  path in  $G-x$



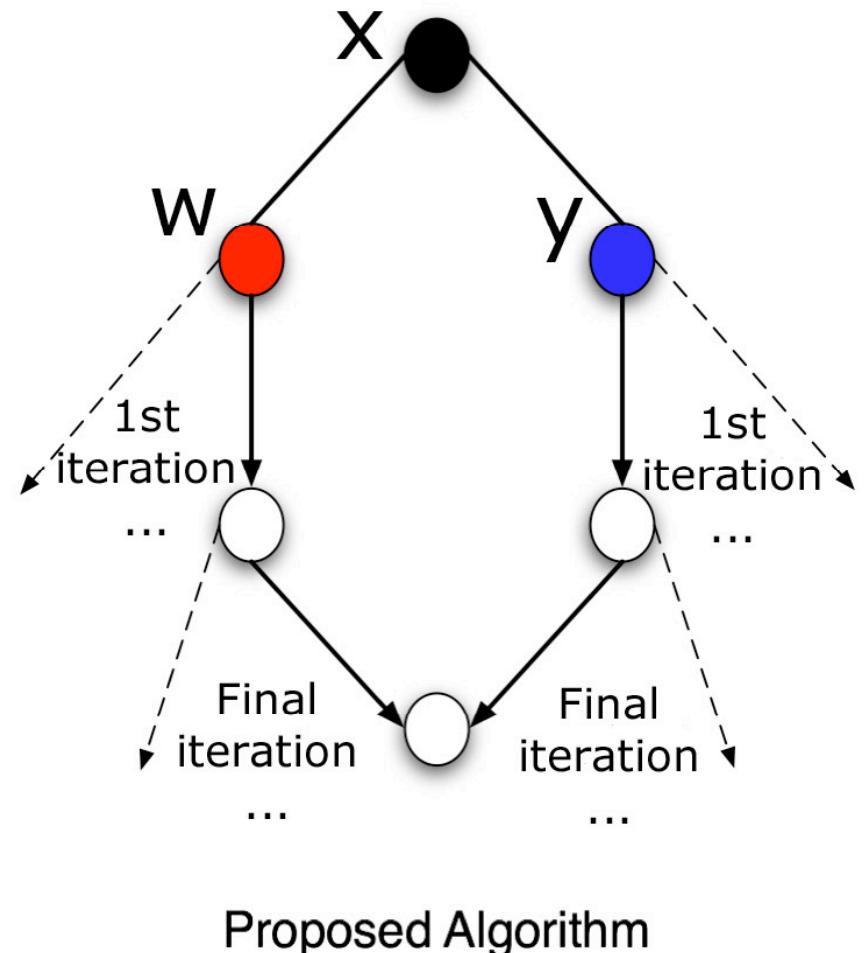
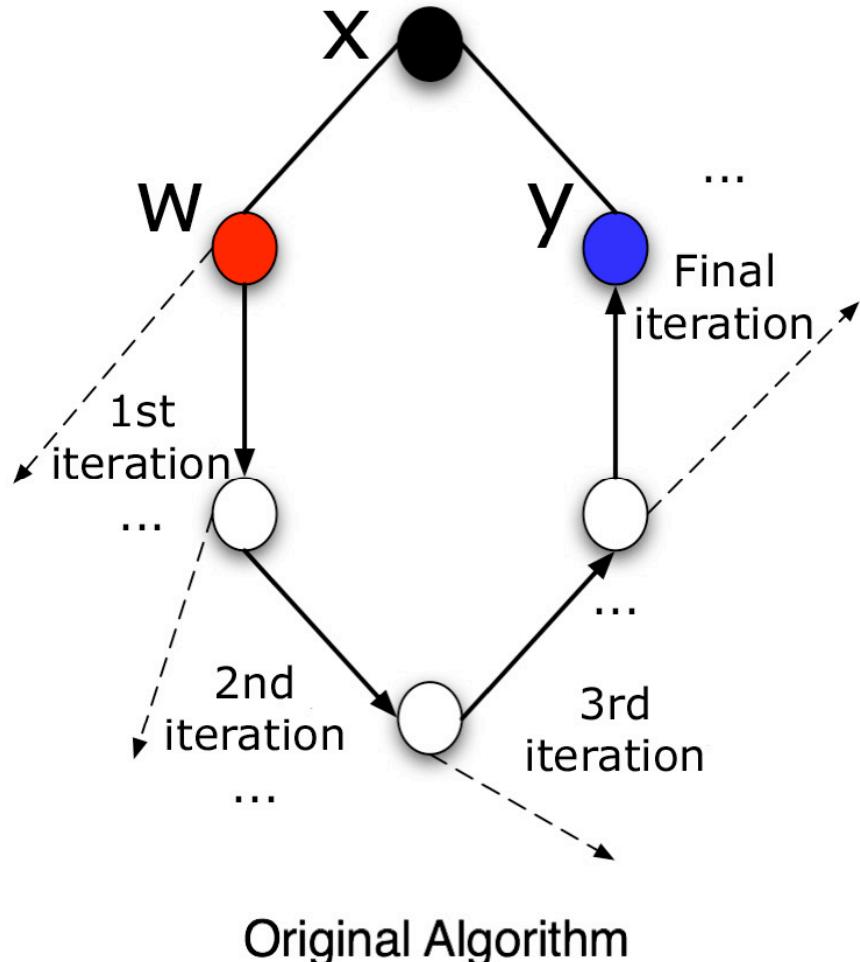
# Ring-based Data Mining

Shortest-path ring analysis of intermediate-range order (IRO) in disordered materials



Correlation between IRO in neutron scattering & ring distribution

# Fast Ring Analysis: Dual-Tree Expansion



# DTE Algorithm

**Algorithm** dual\_tree\_expansion()

**Input:**

$V$  = Set of all vertices (i.e., atoms)  
 $R_c$  = Ring cutoff range (Euclidean)  
 $R_{bc}$  = Bond cutoff distance (Euclidean)  
 $L_{MAX}$  = Maximum length of ring (integer)  
 $P$  = Number of compute nodes

**Output:**

The K-ring statistics for all vertices in the network  
List of atoms with abnormal ring profile

**Variables:**

$\text{Neighbors}(V)$  = Set of vertices that share an edge with vertex  $V$   
 $K_v(p)$  = Number of  $p$ -member rings that go through vertex  $V$   
 $L_b$  = Length of the ring formed with path  $(V_i, V, V_j)$

**Steps:**

- 0 coarse grained spatial decomposition of atoms on  $P$  compute nodes with a thin boundary extension of  $R_c$  distance (This step is for the parallel version only)
- 1 create adjacency list  $G$  for all node in  $V_o$  using  $R_{bc}$  as cutoff distance
- 2 for every vertex  $V \in V_o$ 
  - for each vertex pair  $V_i$  and  $V_j$  in  $\text{Neighbors}(V)$  do
    - $A_1 = \{V_i\}$
    - $A_2 = \{V_j\}$
    - $L_b = 0$
    - while  $(A_1 \cap A_2 = \emptyset \text{ AND } L_b < L_{MAX})$  do
      - $L_b = L_b + 2$
      - if  $(A_1 \cap \text{Neighbors}(A_2) \neq \emptyset \text{ OR } A_2 \cap \text{Neighbors}(A_1) \neq \emptyset)$ 
        - $L_b = L_b + 1$
        - break
      - else if  $(\text{Neighbors}(A_1) \cap \text{Neighbors}(A_2) \neq \emptyset)$ 
        - $L_b = L_b + 2$
        - $A_1 = \text{Neighbors}(A_1)$
        - $A_2 = \text{Neighbors}(A_2)$
        - if  $(L_b < L_{MAX}) \quad ++ K_v(L_b)$

# Spatial Hash-Function Tagging

**Algorithm** spatial hash function tagging (SHAFT)

**Input:**

$C(V)$  = 3D coordinates of all vertices (i.e., atoms)

$R_c$  = Ring cutoff range (Euclidean)

$R_{bc}$  = Bond cutoff distance (Euclidean)

$L_{MAX}$  = Maximum length of ring (integer)

**Output:**

The integer index that is unique for all vertices in the maximum ring span

$$b = R_{lower} / \sqrt{3}$$

$$c = R_{upper} L_{max}$$

$$m = \lceil c/b \rceil$$

**Step:**

for each vertex

    for each spatial dimension  $i$  from 1 to 3

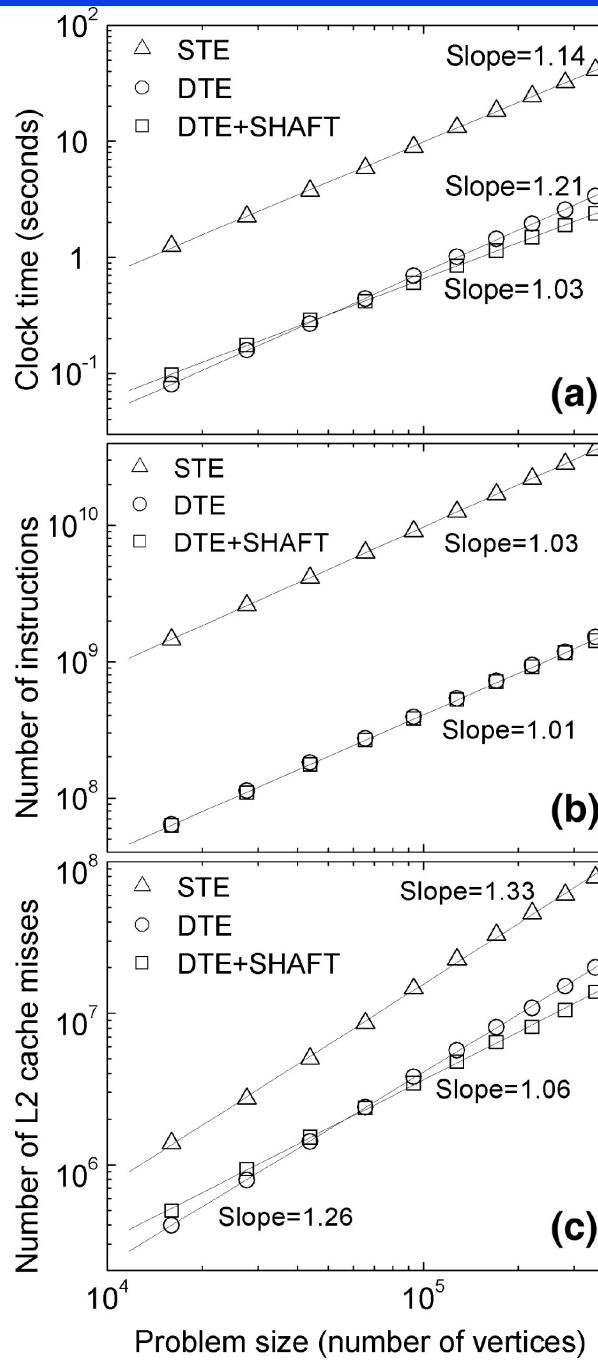
$$q_i = \lfloor C_i/b \rfloor$$

$$q_i \% = m$$

$$\text{return } q = q_3 \times m^2 + q_2 \times m + q_1$$

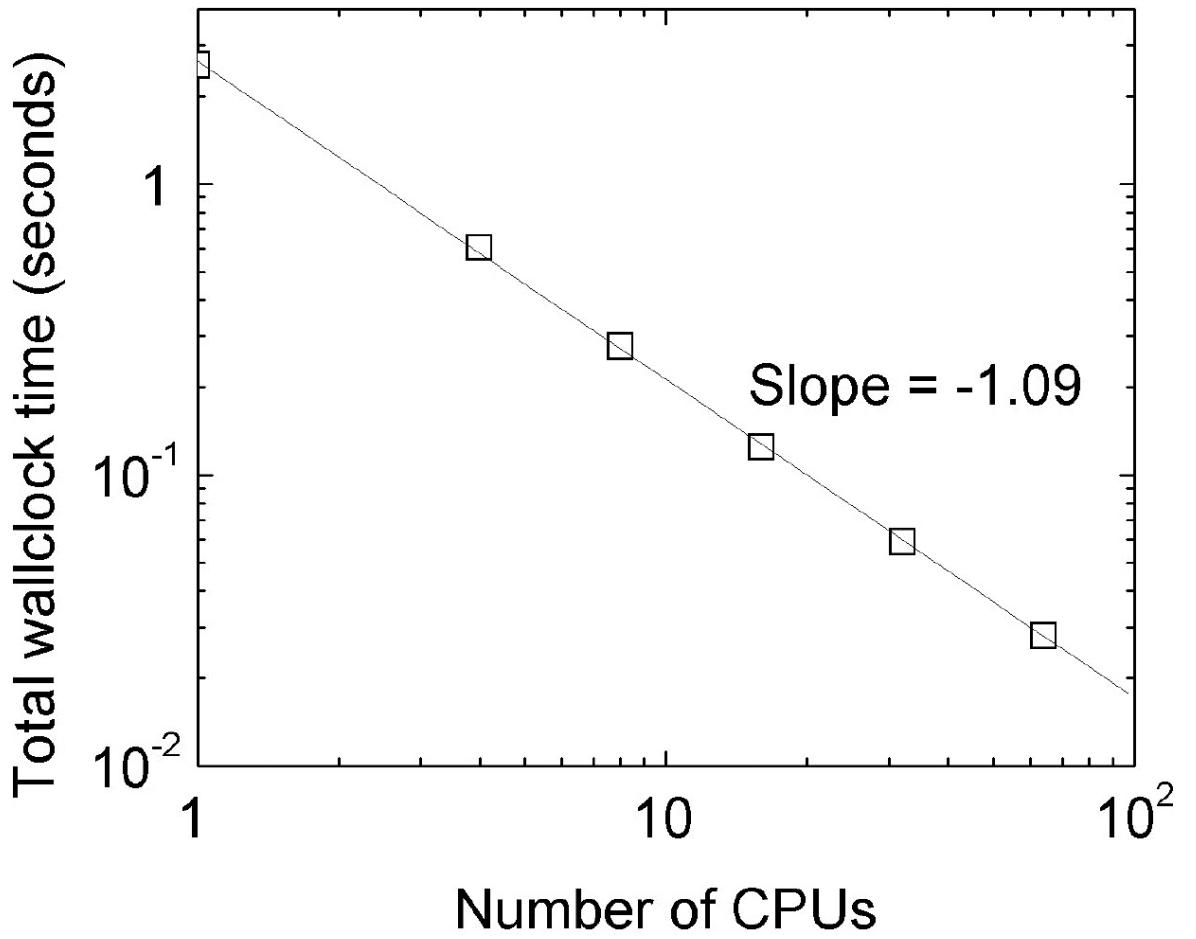
0	1	2	3	4	0	1	2	3	4
5	6	7	8	9	5	6	7	8	9
10	11	12	13	14	10	11	12	13	14
15	16	17	18	19	15	16	17	18	19
20	21	22	23	24	20	21	22	23	24
0	1	2	3	4	0	1	2	3	4
5	6	7	8	9	5	6	7	8	9
10	11	12	13	14	10	11	12	13	14
15	16	17	18	19	15	16	17	18	19
20	21	22	23	24	20	21	22	23	24

# Numerical Tests



Linear scaling  
on the problem size

Superlinear (strong) scaling  
on the number of CPUs

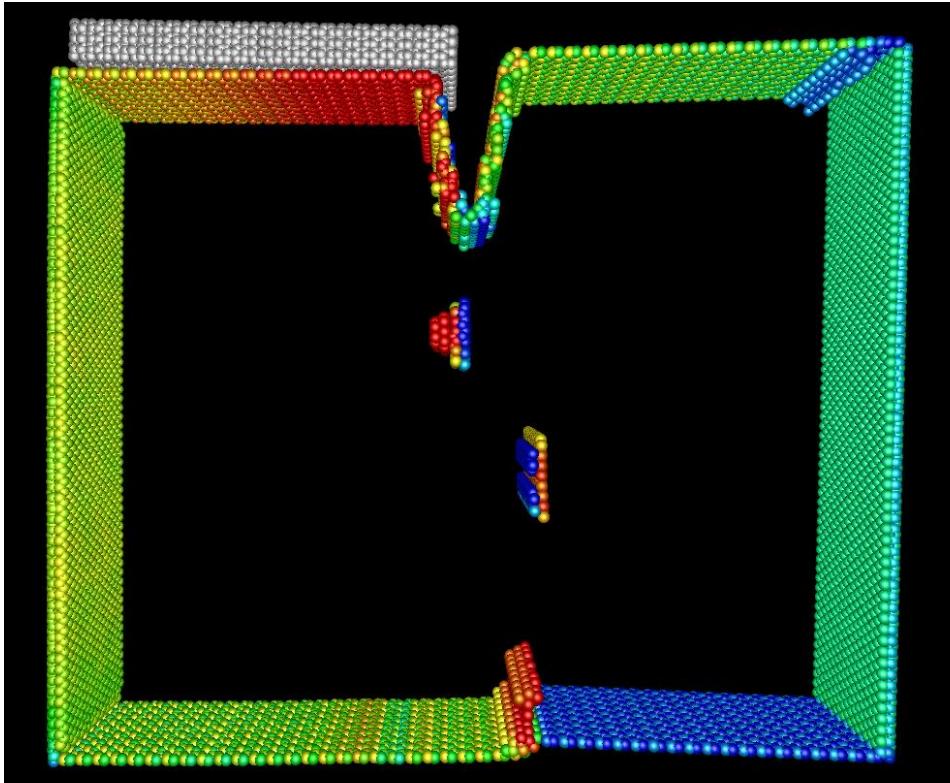


# Dislocation Mining

---

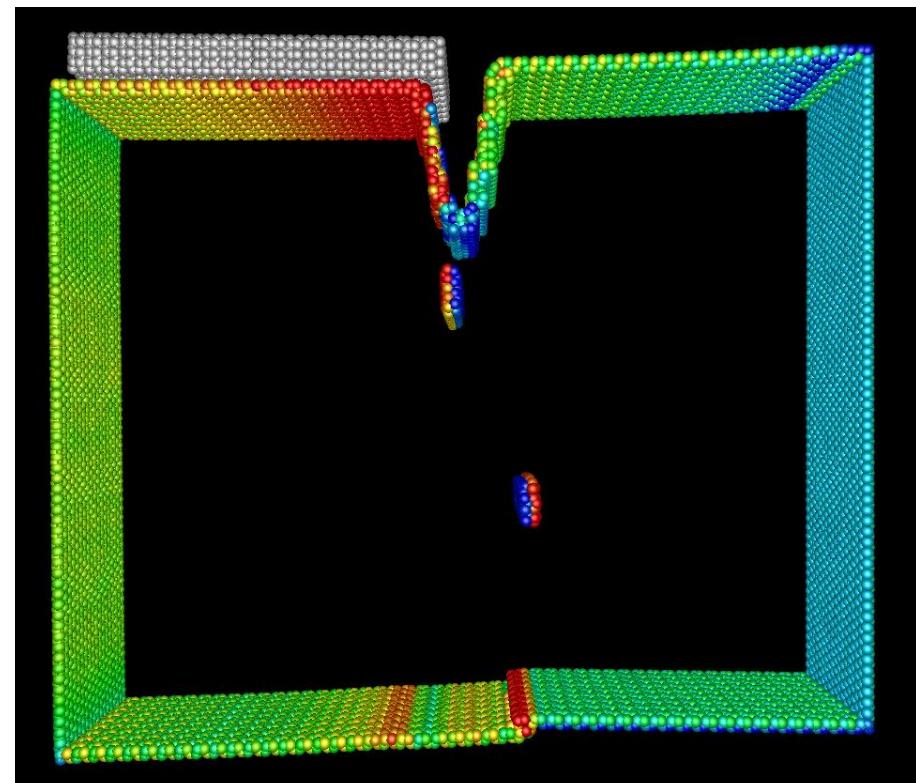
Based on  
potential energy

Shown atoms with high energy compared to bulk



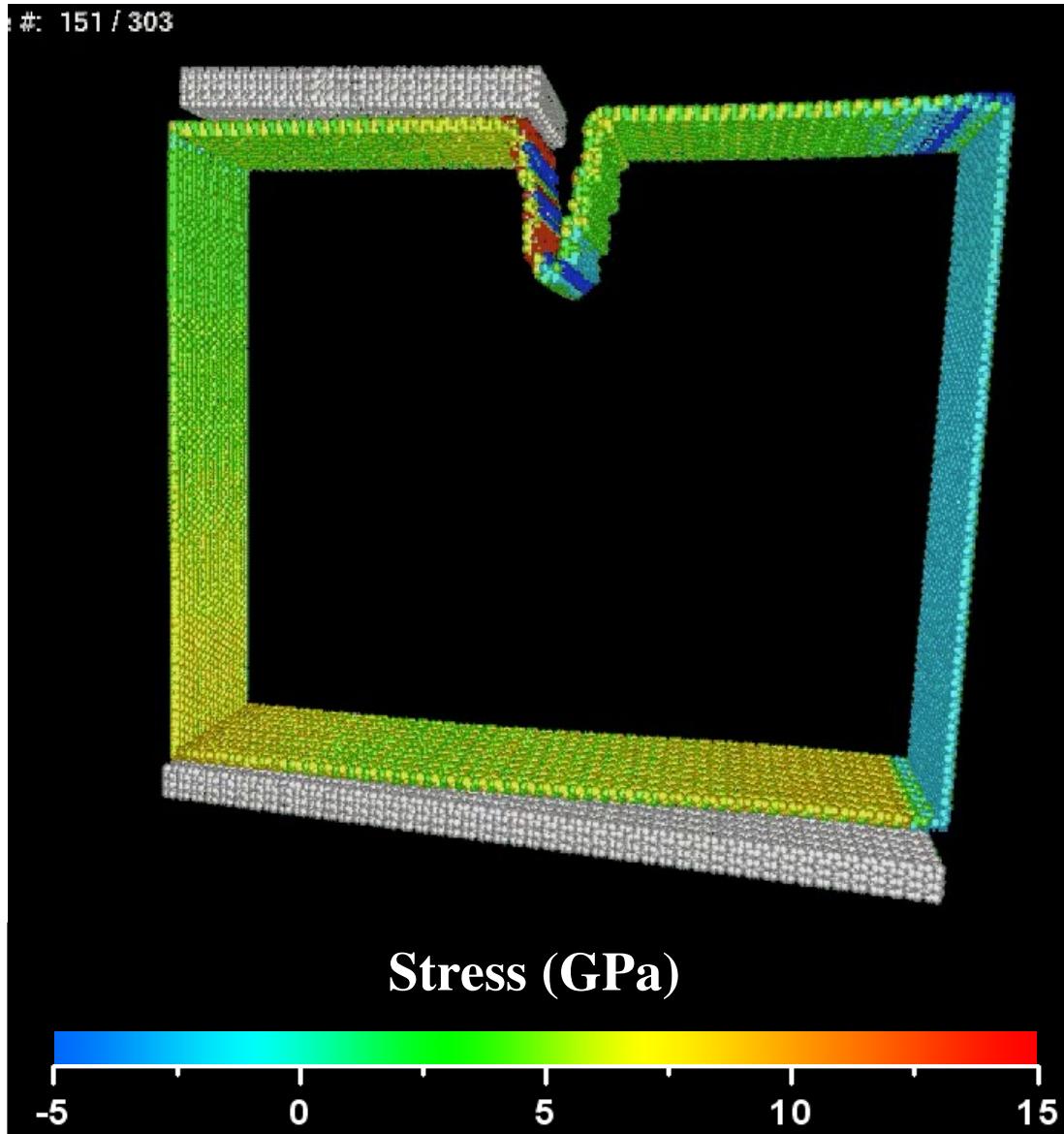
Based on  
shortest-path ring statistics

Shown atoms with less than 12 6-membered rings



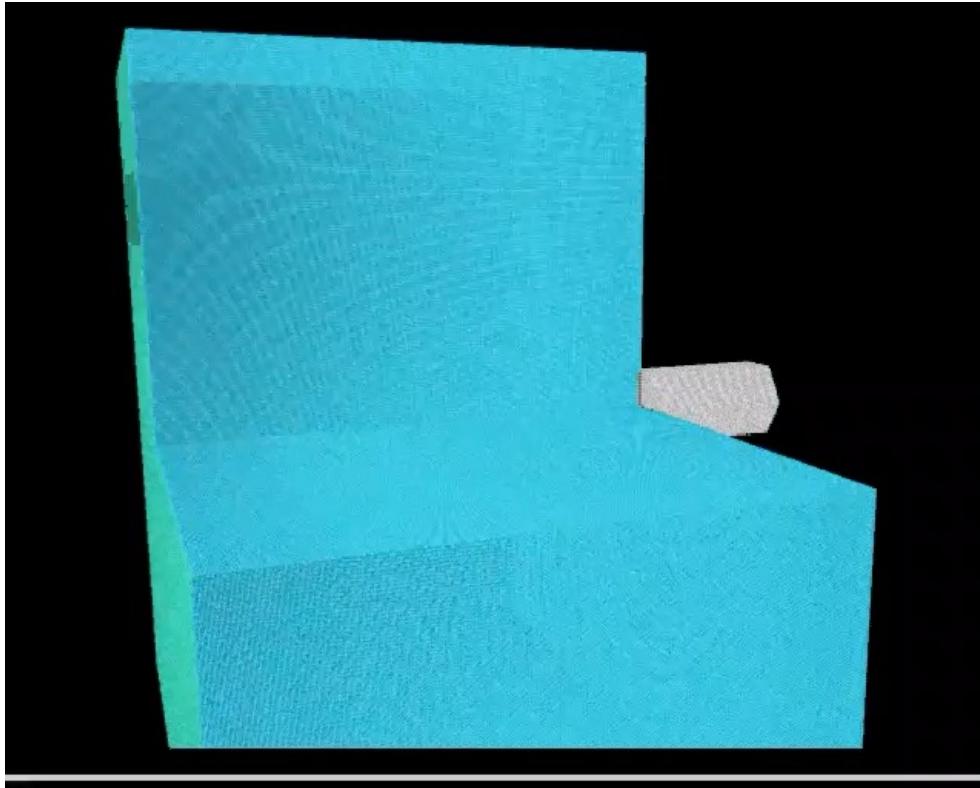
# 100 km/s Impact on Notched AlN

- Dislocation nucleation & emission from notch during impact
- Dislocations & surface atoms mined by ring statistics

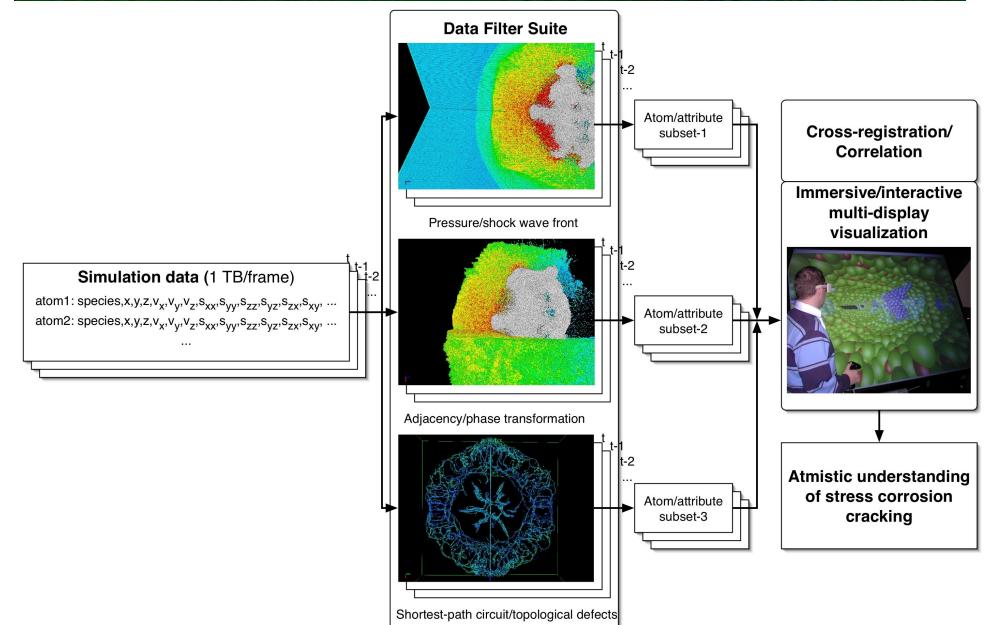
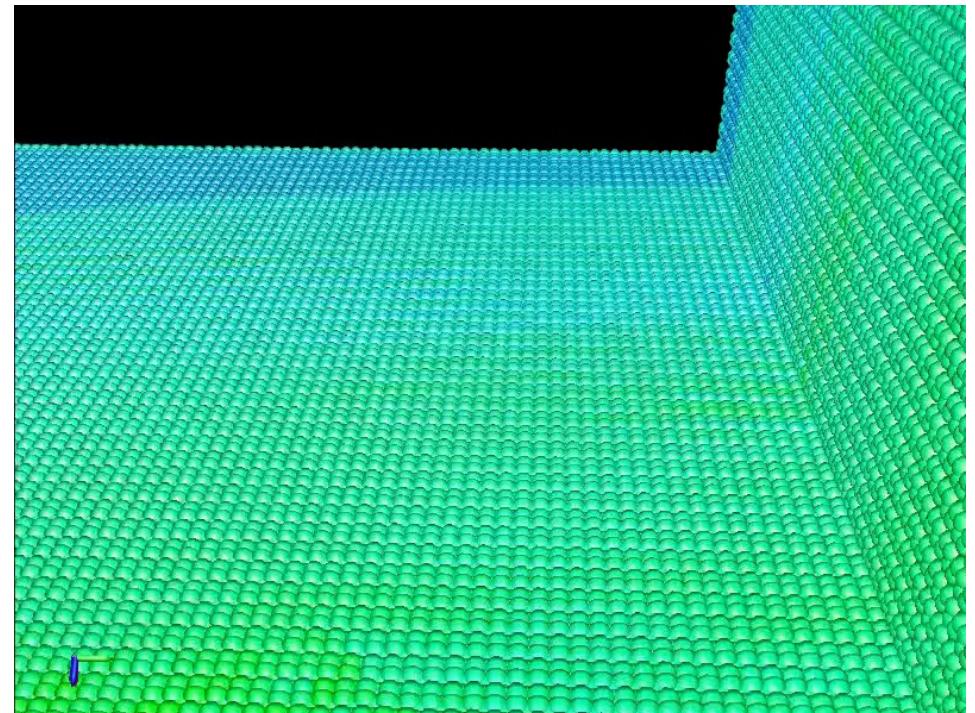


# Impact-Damage Tolerant Ceramics?

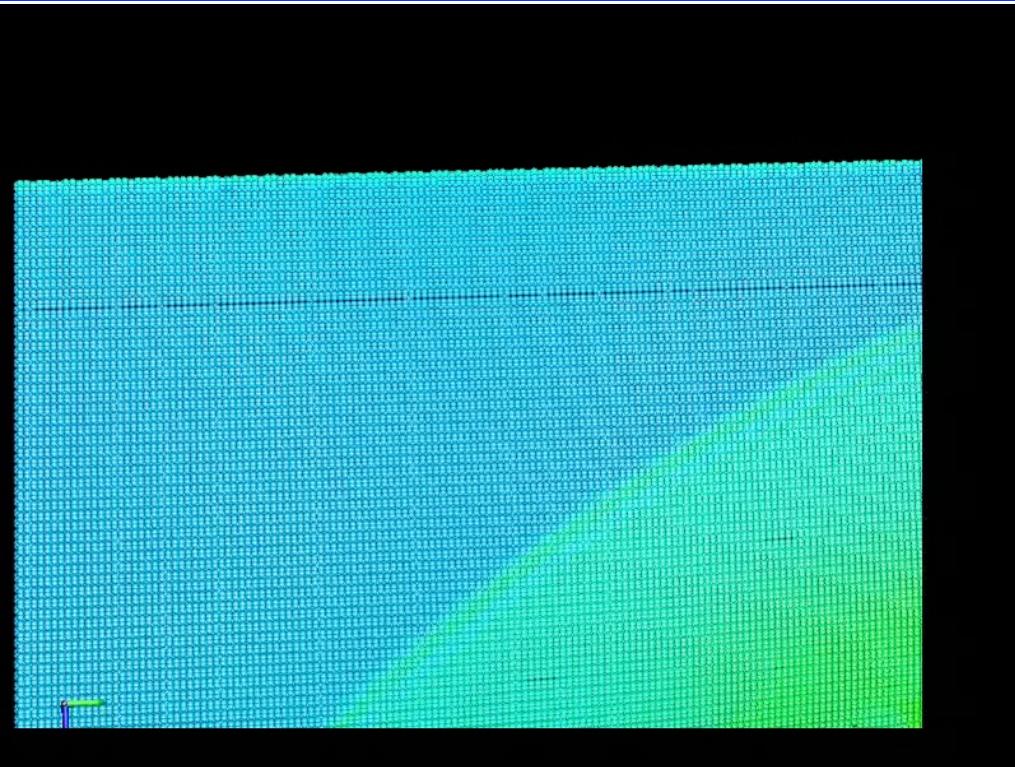
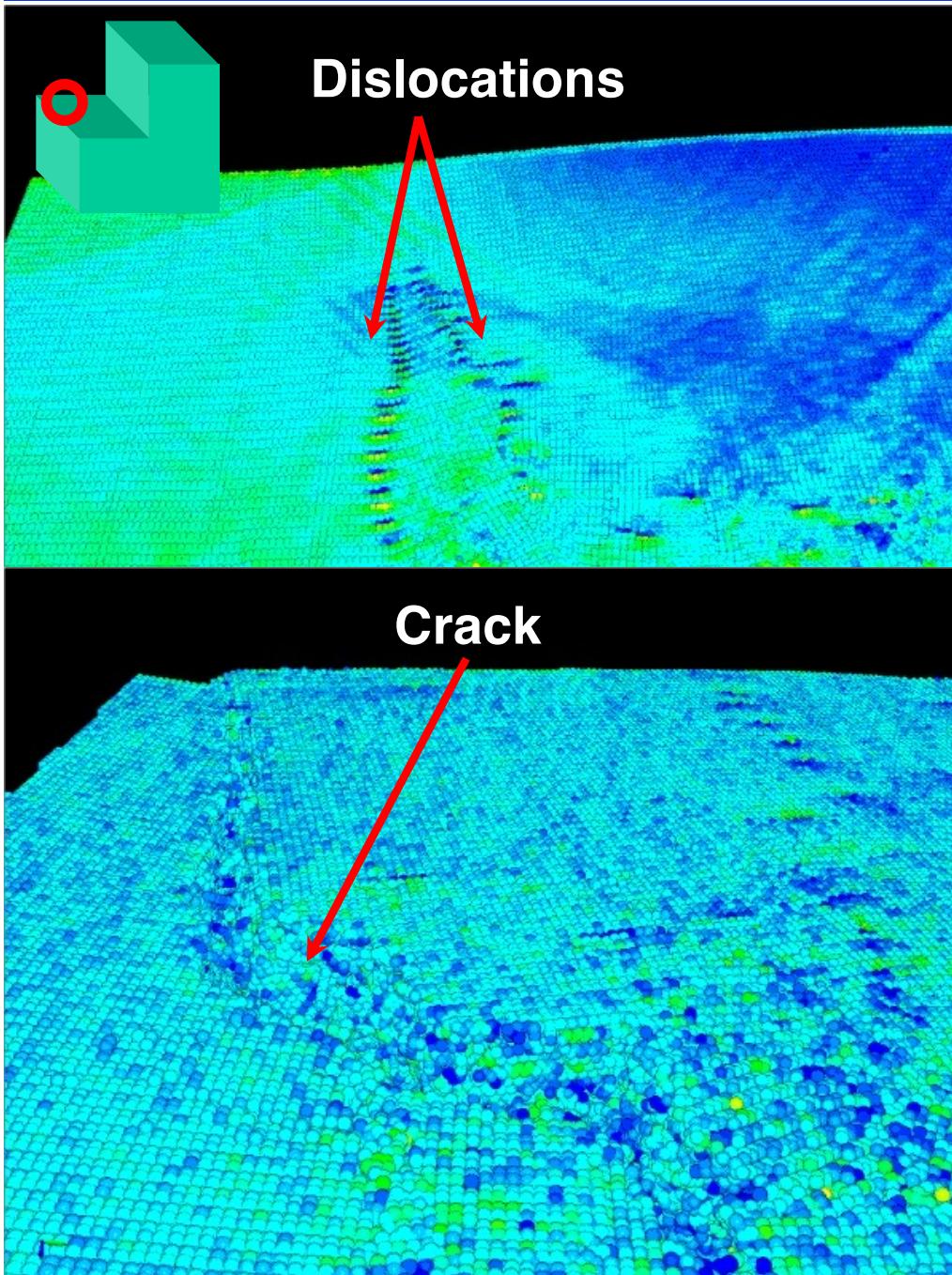
Inverse problem: design materials with desired properties



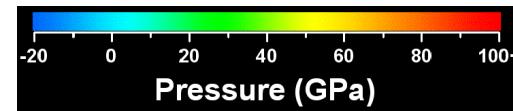
209 million atom MD of hypervelocity impact in AlN for the design of light-weight ceramic armors



# Crack Nucleation at Kink Bands

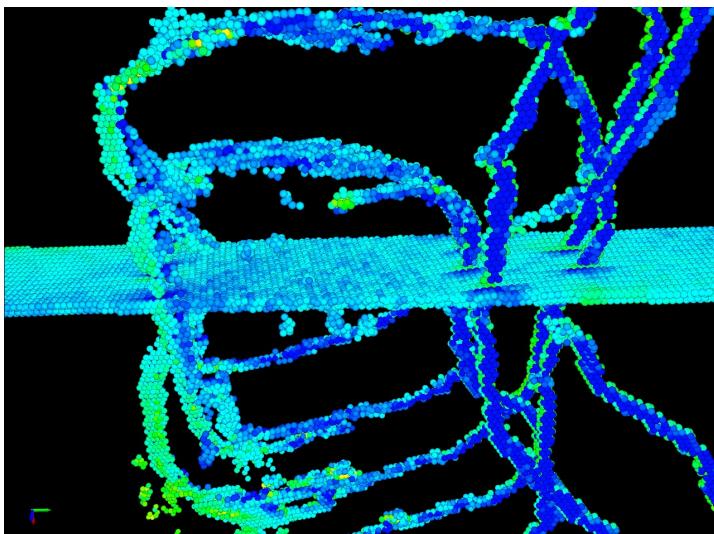
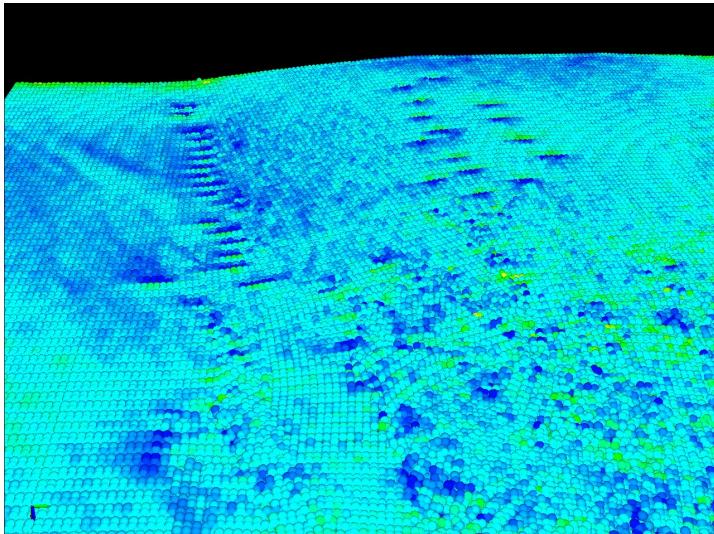


- Series of dislocation dipoles with opposite Burgers vectors form a kink band to releases stress
- Tilt grain boundaries of the kink bands act as sources of mode-II (shear) crack nucleation

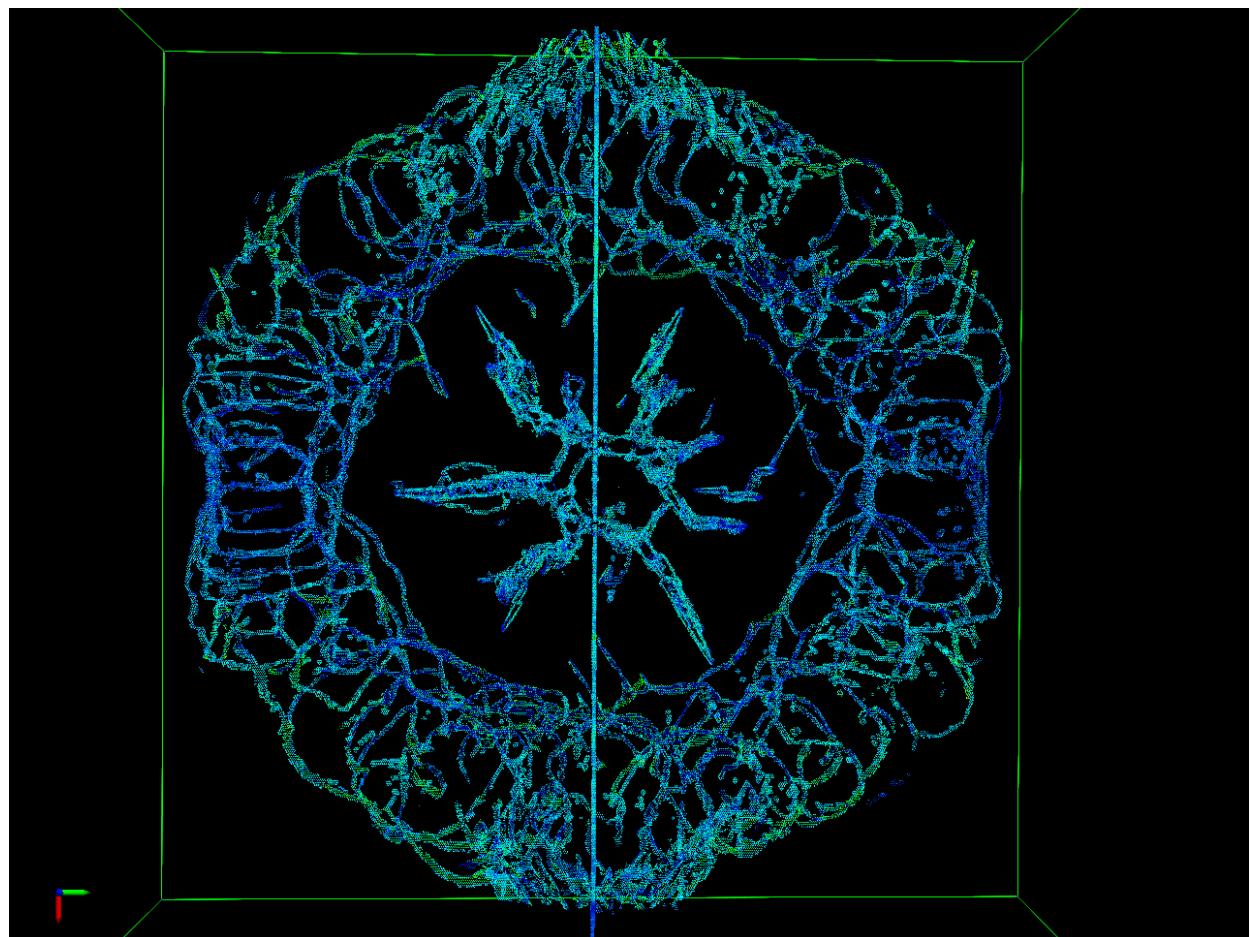


# Dislocation Loops at Kink Bands

Graph (shortest-path circuit) based mining of topological defects



Atoms participating in  
non-6-member circuits



Dislocation network

# Nanoindentation on Nanophase SiC

## Superhardness

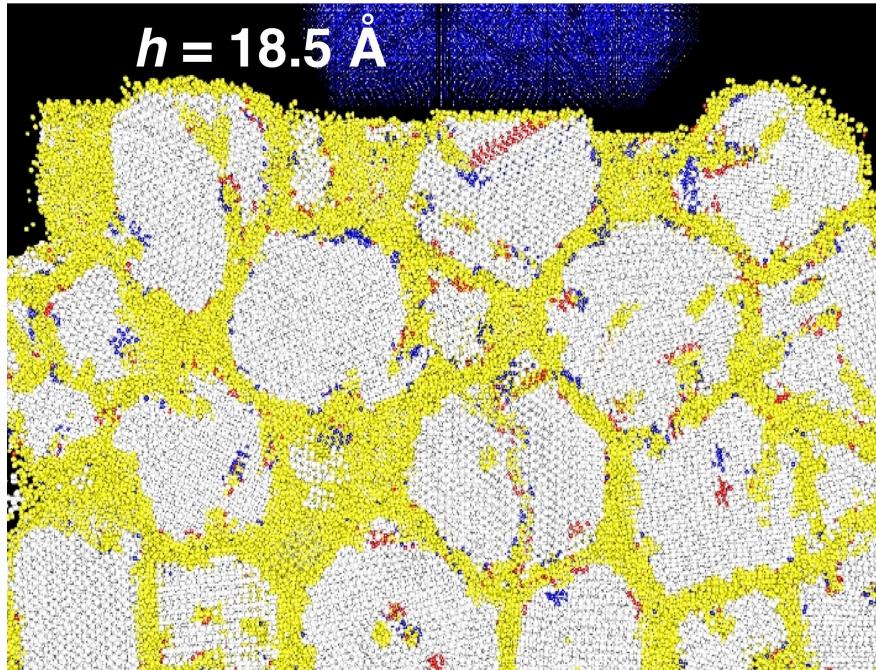
MD: 39 GPa

(grain size,  $d = 8 \text{ nm}$ )

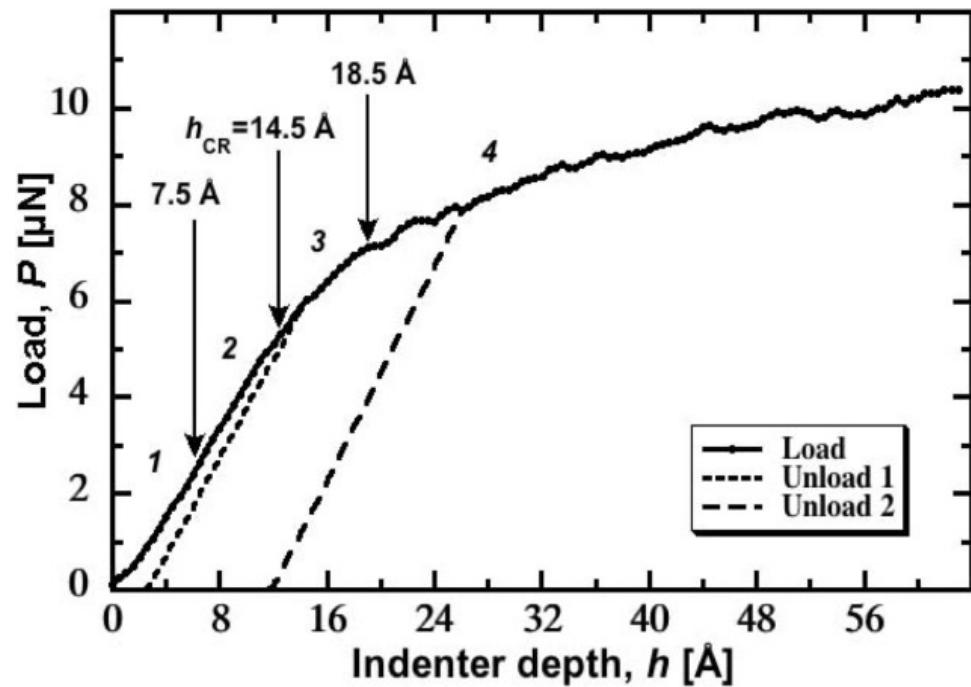
Expt.: 30-50 GPa

( $d = 5\text{-}20 \text{ nm}$ )

[Liao *et al.*, APL, '05]



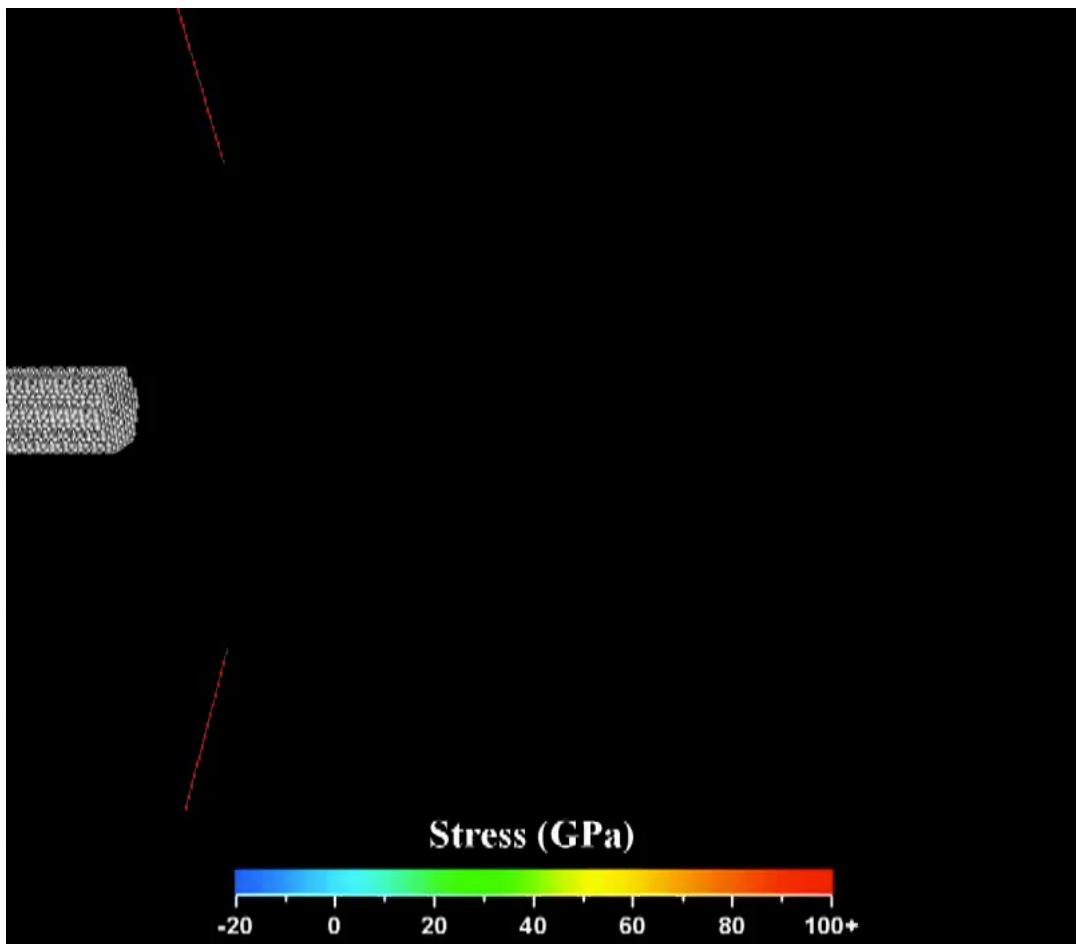
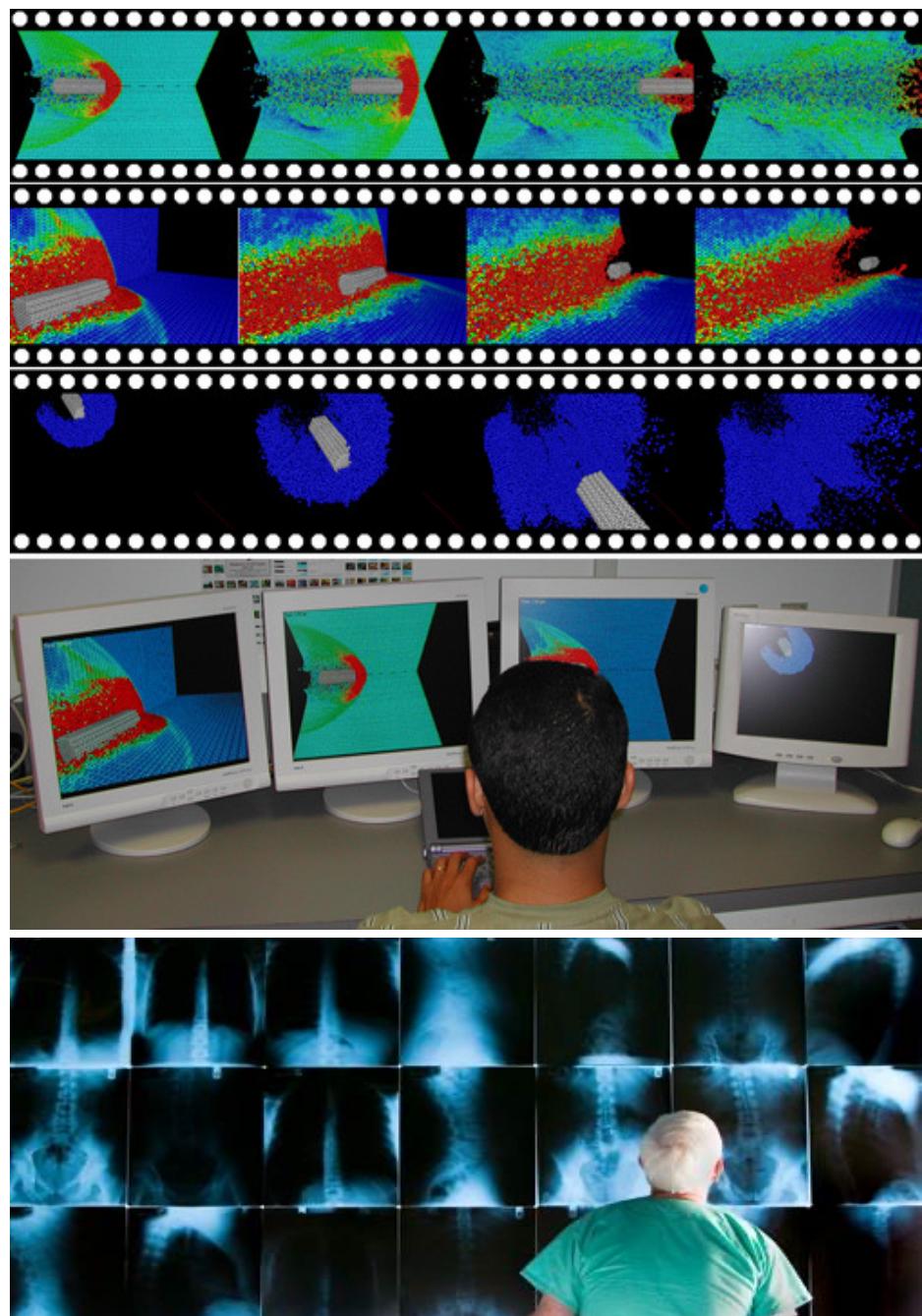
## Load-displacement curve



Crossover from intergrain continuous response to intragranular discrete response

Szlufarska, Nakano & Vashishta, *Science* 309, 911 ('05)

# Multimodal Multidisplay Visualization



Hypervelocity penetration  
through an AlN plate

# Singular Value Decomposition

---

Aiichiro Nakano

*Collaboratory for Advanced Computing & Simulations  
Department of Computer Science  
Department of Physics & Astronomy  
Department of Quantitative & Computational Biology  
University of Southern California*

Email: [anakano@usc.edu](mailto:anakano@usc.edu)

**Goal: Another matrix decomposition (SVD)  
for low-rank matrix approximation**

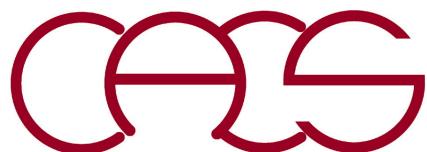
*cf.* Eigen decomposition

$$A = Q \begin{bmatrix} \Sigma \\ 0 \end{bmatrix} Q^T$$

QR decomposition

$$A = Q \begin{bmatrix} R \\ 0 \end{bmatrix}$$

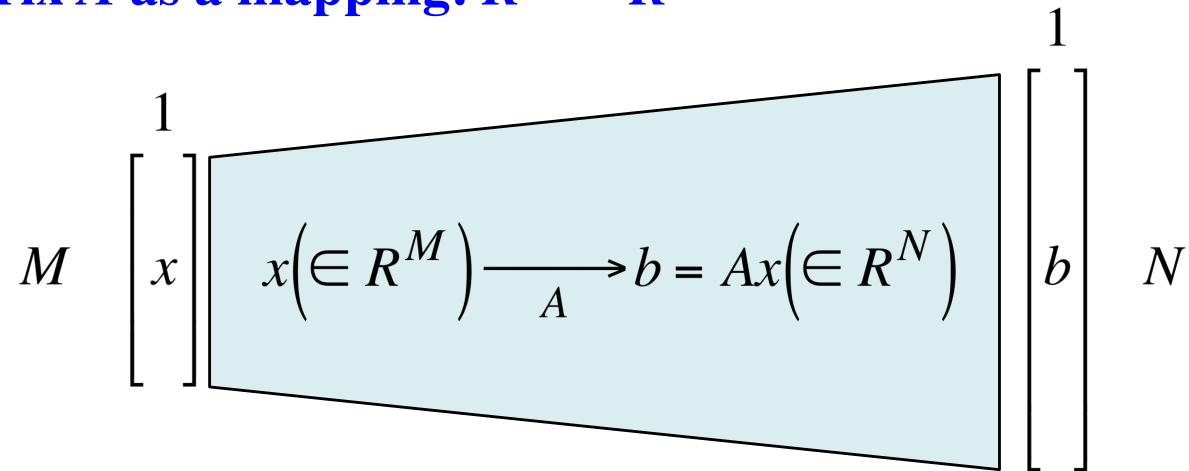
See note on “least square fit” & *Numerical Recipes Sec. 2.6*



# Rank of a Matrix

---

- $N \times M$  matrix  $A$  as a mapping:  $R^M \rightarrow R^N$



- **Range of  $A$ :** Vector space  $\{b = Ax | \forall x\}$
- **Rank of  $A$ :** Number,  $m$ , of linearly-independent vectors in the range, i.e., how many linearly-independent  $N$ -element vectors are there in the range, such that

$$b = A^\top x = \sum_{v=1}^m c_v |v\rangle$$

# Low Rank Approximations of a Matrix

- Rank-1 approximation:  $NM \rightarrow N + M$

$$N \begin{bmatrix} M \\ \psi \end{bmatrix} \cong \underbrace{\begin{bmatrix} u \\ \vdots \end{bmatrix}}_{N \times 1} \underbrace{\begin{bmatrix} v \\ \vdots \end{bmatrix}}_{1 \times M}$$

$|u\rangle\langle v| \forall x \rangle \propto |u\rangle$

- Rank-2 approximation:  $NM \rightarrow 2(N + M)$

$$\begin{bmatrix} \psi \end{bmatrix} \cong \begin{bmatrix} u_1 \\ \vdots \end{bmatrix} w_1 \begin{bmatrix} v_1 \end{bmatrix} + \begin{bmatrix} u_2 \\ \vdots \end{bmatrix} w_2 \begin{bmatrix} v_2 \end{bmatrix}$$

- Rank- $m$  ( $m \ll N, M$ ) approximation:  $NM \rightarrow m(N + M)$

$$\begin{bmatrix} \psi \end{bmatrix} \cong \sum_{v=1}^m \begin{bmatrix} u_v \\ \vdots \end{bmatrix} w_v \begin{bmatrix} v_v \end{bmatrix}$$

# Singular Value Decomposition

- **Problem:** Optimal approximation of an  $N \times M$  matrix  $\psi$  of rank- $m$  ( $m \ll N$ )?
- **Theorem:** An  $N \times M$  matrix  $\psi$  (assume  $N \geq M$ ) can be decomposed as

$$\psi = UDV^T = \sum_{\nu=1}^M U_{i\nu} d_\nu V_{j\nu} = \sum_{\nu=1}^M u_i^{(\nu)} d_\nu v_j^{(\nu)}$$

where  $U \in \mathbf{R}^N \times \mathbf{R}^M$  &  $V \in \mathbf{R}^M \times \mathbf{R}^M$  are column orthogonal &  $D$  is diagonal

$$N \begin{bmatrix} \psi \end{bmatrix}_{N \times M} = U \begin{bmatrix} d_1 & & \\ & \ddots & \\ & & d_M \end{bmatrix}_{M \times M} V^T \quad \text{See appendix on polar & singular decompositions}$$

- **Theorem:** Sort the SVD diagonal elements in descending order,  $d_1 \geq d_2 \geq \dots \geq d_M \geq 0$ , & retain the first  $m$  terms

$$\psi^{(m)} \equiv \sum_{\nu=1}^m u_i^{(\nu)} d_\nu v_j^{(\nu)T}$$

which is optimal among  $\forall$  rank- $m$  matrices in the 2-norm sense with the error

$$\min_{\text{rank}(A)=m} \|A - \psi\|_2 = \|\psi^{(m)} - \psi\|_2 = d_{m+1}$$

cf. [singular.c](#) & [svdcmp.c](#)

[cc -o singular singular.c svdcmp.c -lm](#)

**Use the program!**

# SVD for Image Compression



Original Image



5 Iterations



10 Iterations

**D. Richards & A. Abrahamsen**



20 Iterations

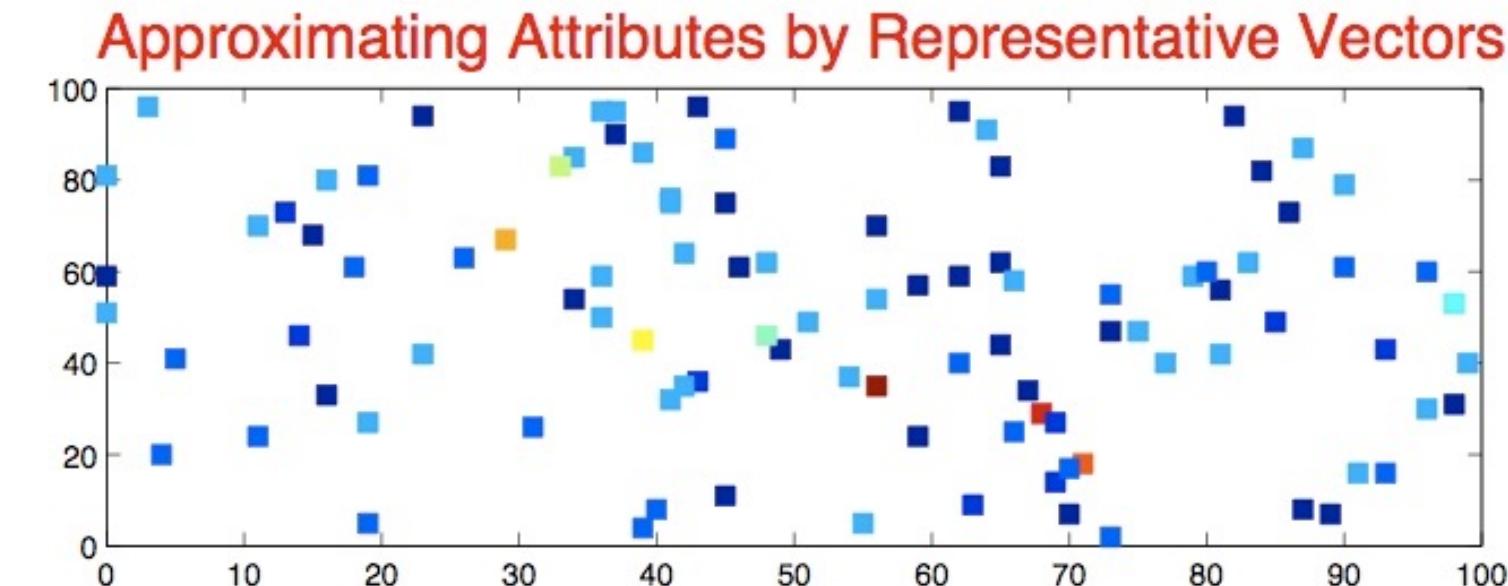
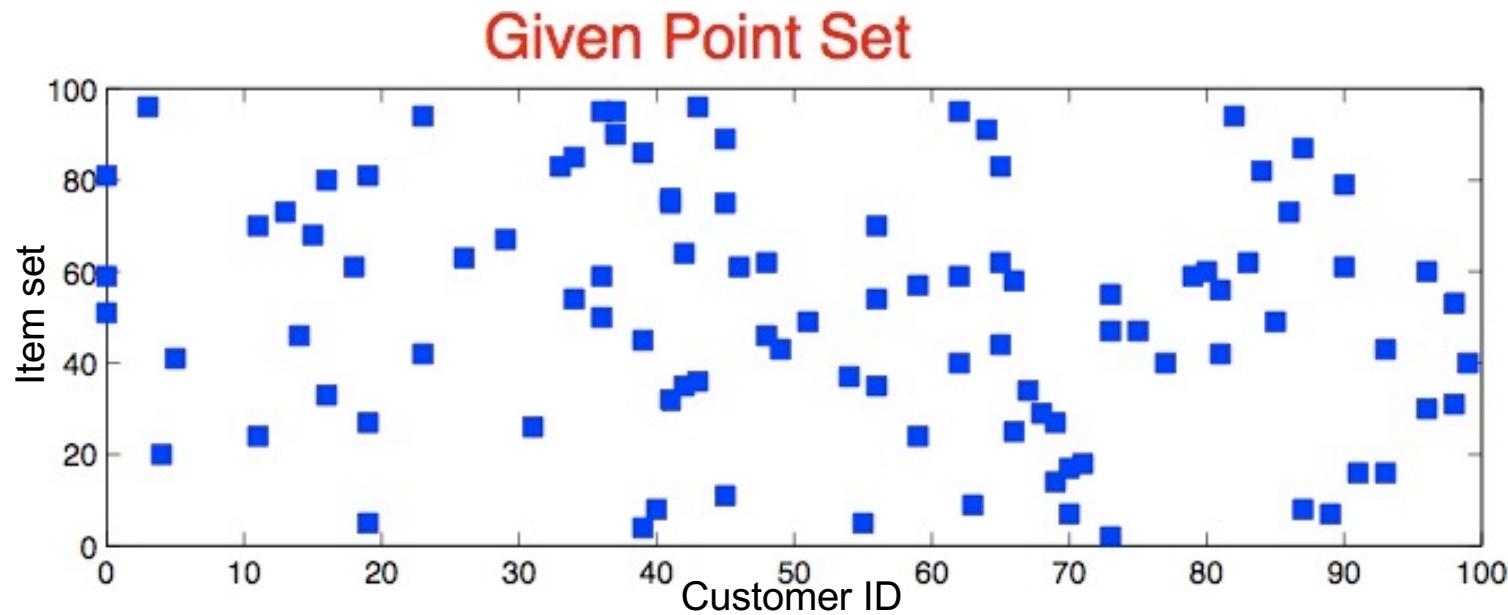


60 Iterations



100 Iterations

# SVD in Data Mining



# Reduced Density Matrix

- Quantum system coupled to an environment



$$\{|i\rangle = \psi_i(x)|i=1, \dots, N\} \quad \{|j\rangle = \phi_j(X)|j=1, \dots, M\}$$

- Quantum state of block + environment

$$|\psi\rangle = \sum_{i=1}^N \sum_{j=1}^M \psi_{ij} |i\rangle |j\rangle \quad \text{or} \quad \Psi(x, X) = \sum_{i=1}^N \sum_{j=1}^M \psi_{ij} \psi_i(x) \phi_j(X)$$

- Reduced density matrix

$$\begin{aligned} \langle \forall A \rangle &= \sum_i \sum_j \psi_{ij}^* \langle j | A \sum_{i'} \sum_{j'} \psi_{i'j'} | i' \rangle | j' \rangle \\ \text{Arbitrary operator in the block} &= \sum_i \sum_j \sum_{i'} \sum_{j'} \psi_{i'j'} \psi_{ij}^* \langle i | A | i' \rangle \langle j | j' \rangle \delta_{jj'} \\ &= \sum_i \sum_{i'} \sum_j \psi_{i'j} \psi_{ij}^* \langle i | A | i' \rangle \equiv \sum_i \sum_{i'} \rho_{i'i} A_{ii'} = \text{tr}_B(\rho A) \end{aligned}$$

$$\rho_{i'i} \equiv \sum_j \psi_{i'j} \psi_{ij}^* \quad A_{ii'} \equiv \langle i | A | i' \rangle$$

# Low-Rank Approx. to Reduced Density Matrix

---

$$\begin{aligned}\psi \cong \psi^{(m)} &= \sum_{\nu=1}^m u^{(\nu)} d_\nu v^{(\nu)T} & \psi_{ij}^{(m)} &= \sum_{\nu=1}^m u_i^{(\nu)} d_\nu v_j^{(\nu)} \\ \rho = \psi\psi^T \cong \psi^{(m)}\psi^{(m)T} &= \sum_{\nu=1}^m \sum_{\nu'=1}^m u^{(\nu)} d_\nu \left( v^{(\nu)T} v^{(\nu')} \right) d_{\nu'} u^{(\nu')T} \\ &= \sum_{\nu=1}^m \sum_{\nu'=1}^m u^{(\nu)} d_\nu (\delta_{\nu\nu'}) d_{\nu'} u^{(\nu')T} & \sum_{\nu=1}^m u^{(\nu)} d_\nu^2 u^{(\nu)T} &\equiv \rho^{(m)} \\ \rho_{ii'}^{(m)} &= \sum_{\nu=1}^m u_i^{(\nu)} d_\nu^2 u_{i'}^{(\nu)}\end{aligned}$$

- **Density matrix renormalization group** = systematic procedure to accurately obtain a quantum ground state:
  1. Incrementally add environment to a block
  2. Solve the global (= block + environment) ground state
  3. Construct a low-rank approx. to represent the block with reduced d.o.f.

S. R. White, Phys. Rev. B **48**, 10345 ('93);

G. K.-L. Chan & S. Sharma, Annu. Rev. Phys. Chem. **62**, 465 ('11) (Eqs. 3-10)

**More SVD in quantum chemistry**

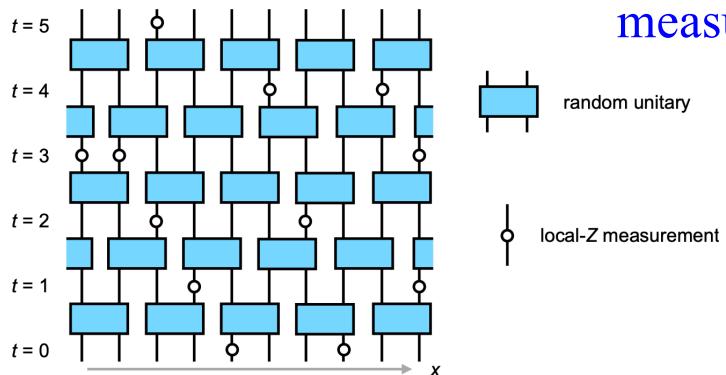
J. F. Gonthier & M. Head-Gordon, J. Chem. Phys. **147**, 144110 ('17)

# Entanglement Entropy & Quantum Computing

- **Entanglement entropy:** A measure of quantum entanglement between two subsystems. If a state describing two subsystems A and B is a *separable* state  $|\Psi_{AB}\rangle = |\phi_A\rangle|\phi_B\rangle$  (i.e., tensor product), then the reduced density matrix  $\rho_A = \text{Tr}_B |\Psi_{AB}\rangle\langle\Psi_{AB}| = |\phi_A\rangle\langle\phi_A|$  is a *pure state*. Thus, the entropy of the state is zero. A reduced density matrix having a non-zero entropy is therefore a signal of the existence of entanglement in the system.
- **Area law:** A quantum state satisfies an *area law* if the leading term of the entanglement entropy grows at most proportionally with the *boundary* between the two partitions. Area laws are remarkably common for ground states of local *gapped* quantum many-body systems. *It greatly reduces the complexity of quantum many-body systems. The density matrix renormalization group and matrix product states, for example, implicitly rely on such area laws.* [https://en.wikipedia.org/wiki/Entropy\\_of\\_entanglement](https://en.wikipedia.org/wiki/Entropy_of_entanglement)

## Measurement-driven entanglement transition in hybrid quantum circuits

“With increasing measurement rate, the volume law phase is unstable to a disentangled area law phase, passing through a single entanglement transition at a critical rate of measurement.”



[Y. Li et al., Phys. Rev. B 100, 134306 \('19\)](#)

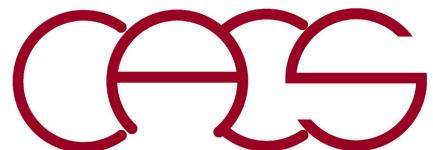
# Machine Learning in Simulation

---

Aiichiro Nakano

*Collaboratory for Advanced Computing & Simulations  
Department of Computer Science  
Department of Physics & Astronomy  
Department of Quantitative & Computational Biology  
University of Southern California*

Email: [anakano@usc.edu](mailto:anakano@usc.edu)

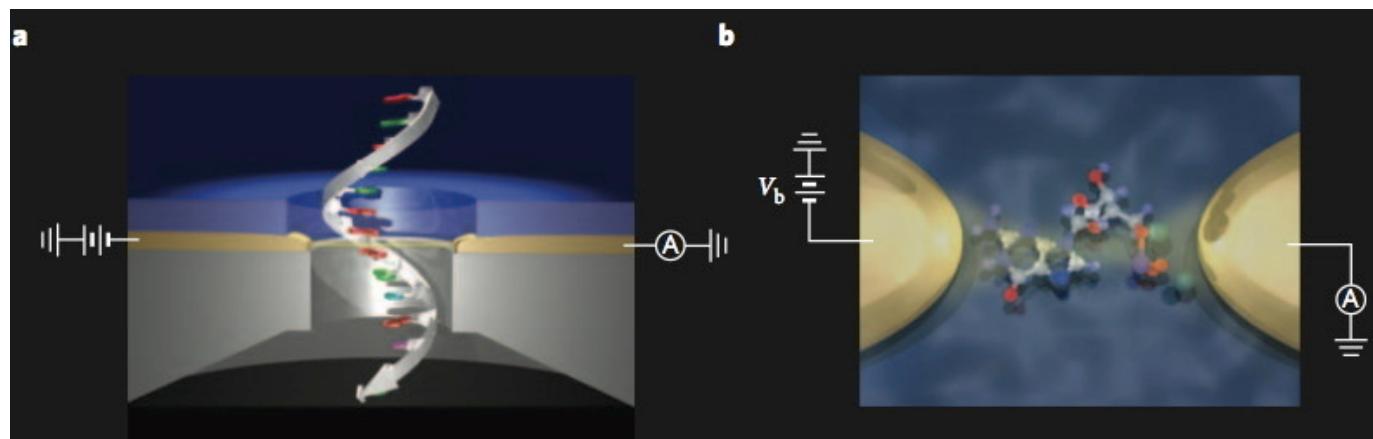


# SVD for Rapid Genome Sequencing

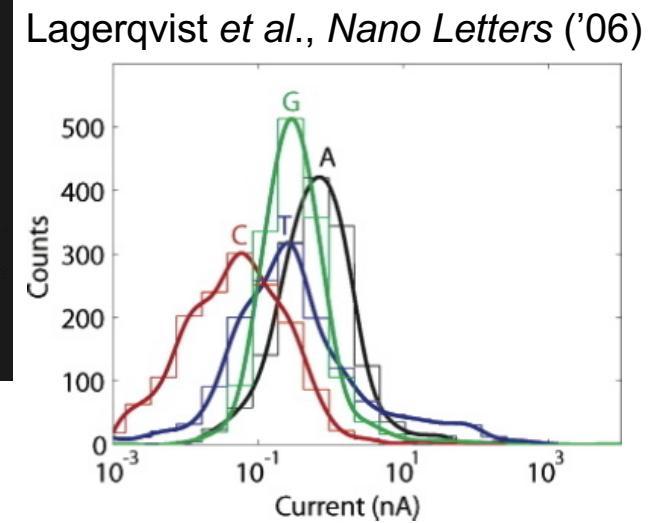
- \$10M Archon X prize for decoding 100 human genomes in 10 days & \$10K per genome (<http://genomics.xprize.org>): Preemptive attack on diseases



- Quantum tunneling current for rapid DNA sequencing?



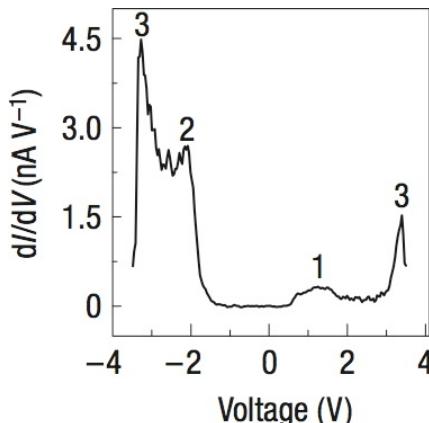
Tsutsui et al., *Nature Nanotechnology* ('10)



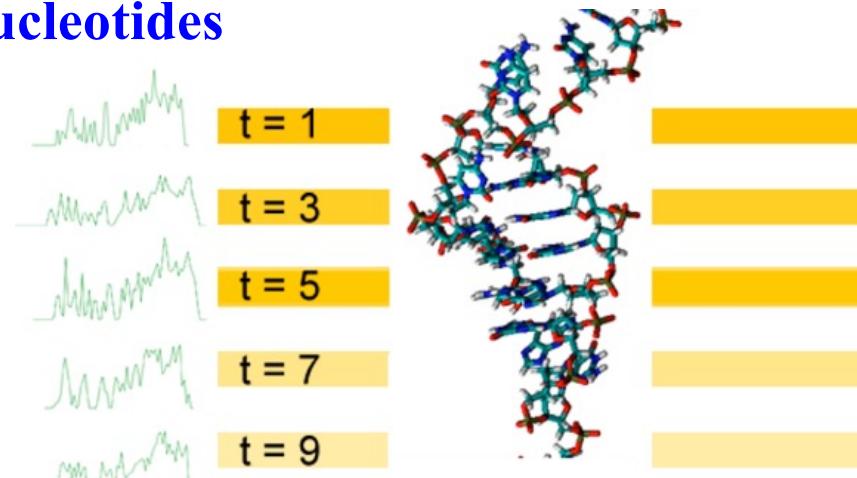
- Tunneling current alone cannot distinguish the 4 nucleotides (A, C, G, T)

# Rapid DNA Sequencing *via* Data Mining

- Use tunneling current ( $I$ )-voltage ( $V$ ) characteristic (or electronic density-of-states) as the ‘fingerprints’ of the 4 nucleotides

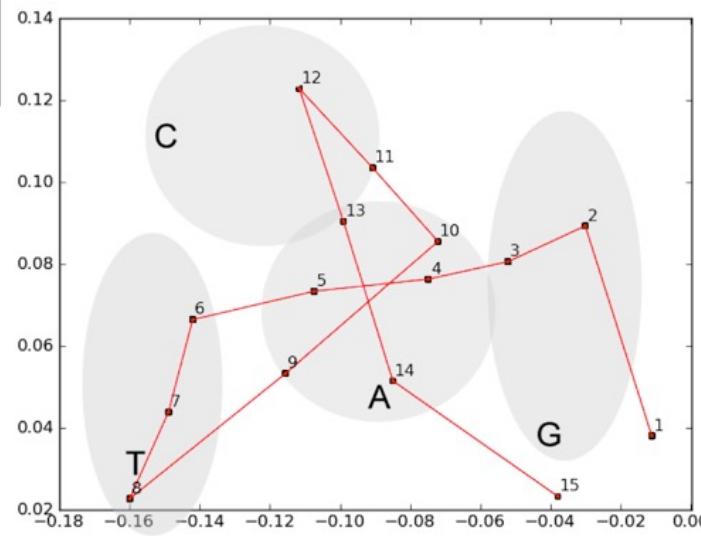
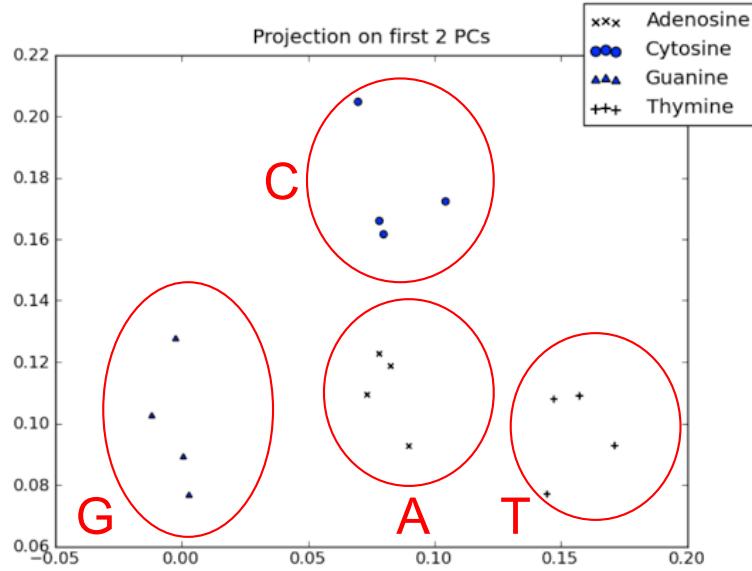


Shapir et al.,  
*Nature Materials* ('08)



- Principal component analysis (PCA) & fuzzy c-means clustering clearly distinguish the 4 nucleotides

H. Yuen et al., *IJCS* 4, 352 ('10)



<http://www.henryyuen.net/>

- Viterbi algorithm for even higher-accuracy sequencing

See [Henry's landmark discovery](#)

# SVD vs. PCA (in Economics)

- SVD of  $N$  (number of companies)  $\times T$  (number of time points) of stock-price time series

$$\underset{T \times N}{[\Sigma]}^T = \underset{T \times N}{[U]} \sum_{N \times N} \underset{N \times N}{[V]}^T$$

- Stock correlation matrix

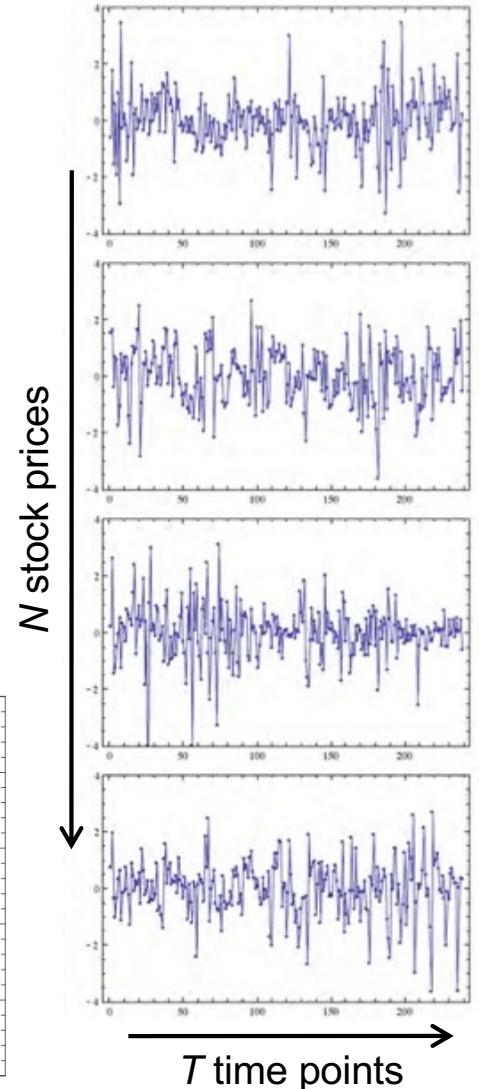
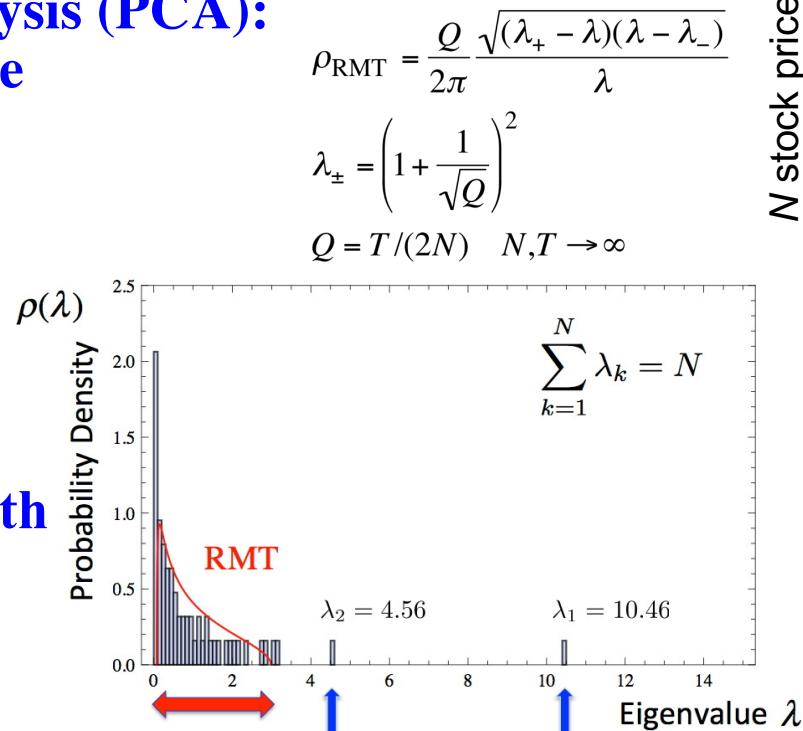
$$\underset{N \times N}{[C]} = \underset{N \times T}{[\Sigma]} \underset{T \times N}{[\Sigma]}^T$$

*Apply it in your area!*

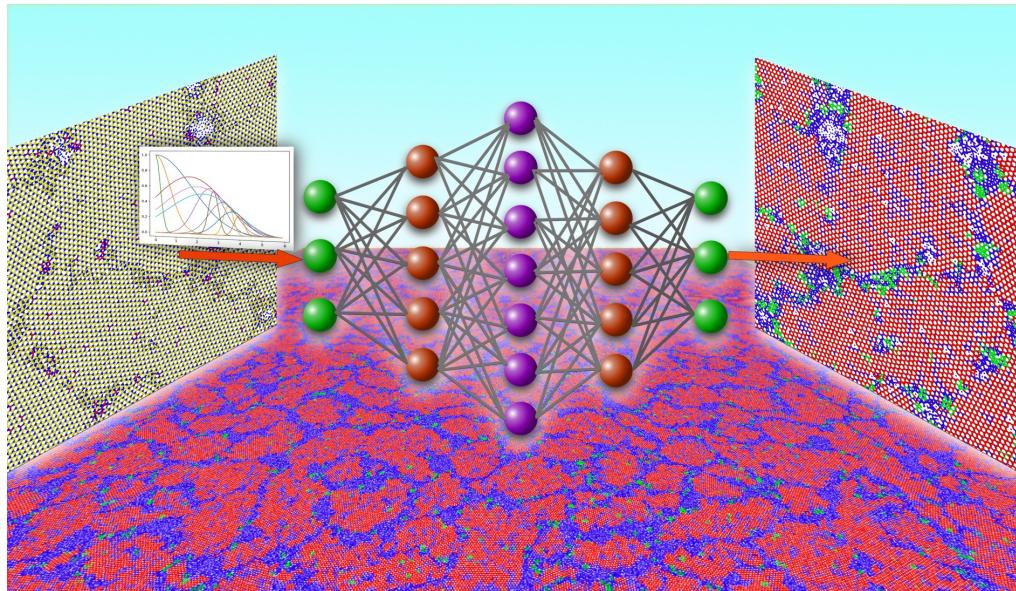
- Principal component analysis (PCA): Eigen decomposition of the correlation matrix

$$\begin{aligned} C &= \underset{I}{\Sigma} \Sigma^T \\ &= V \Sigma \widetilde{U^T U} \Sigma V^T \\ &= V \Sigma^2 V^T \end{aligned}$$

- Compare the spectrum with that of random matrix theory (RMT) for judging statistical significance



# Learning Materials Phases & Defects

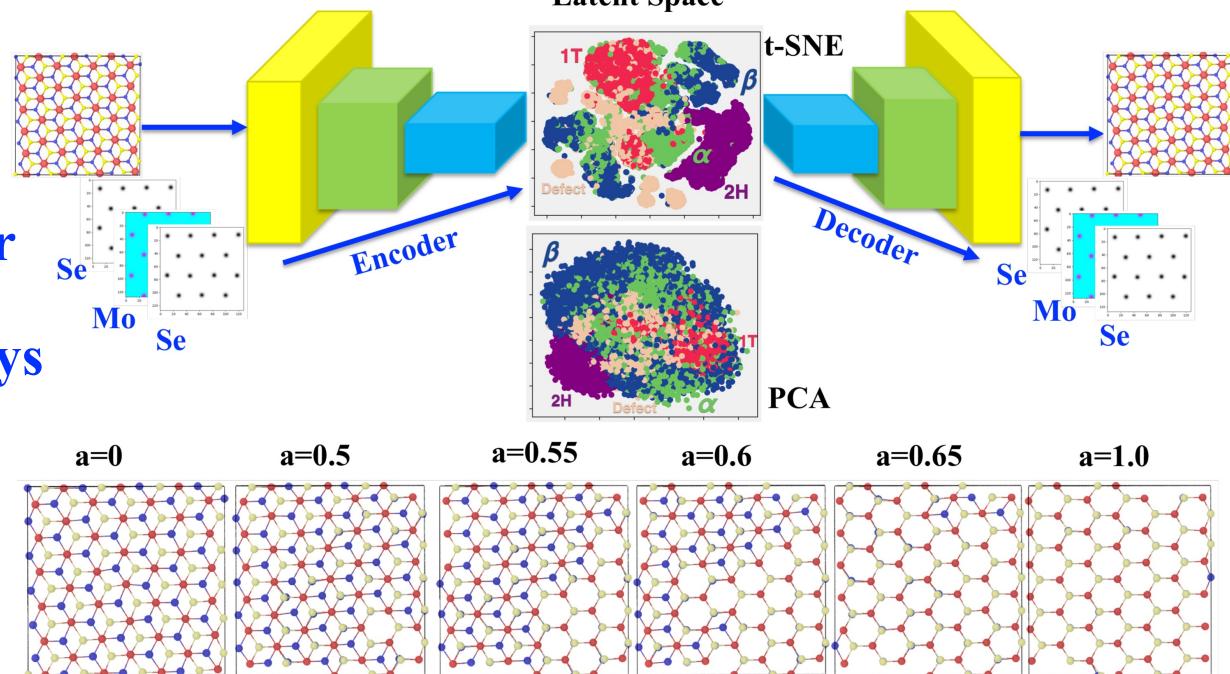


- Feedforward neural network to learn phases from local symmetry functions

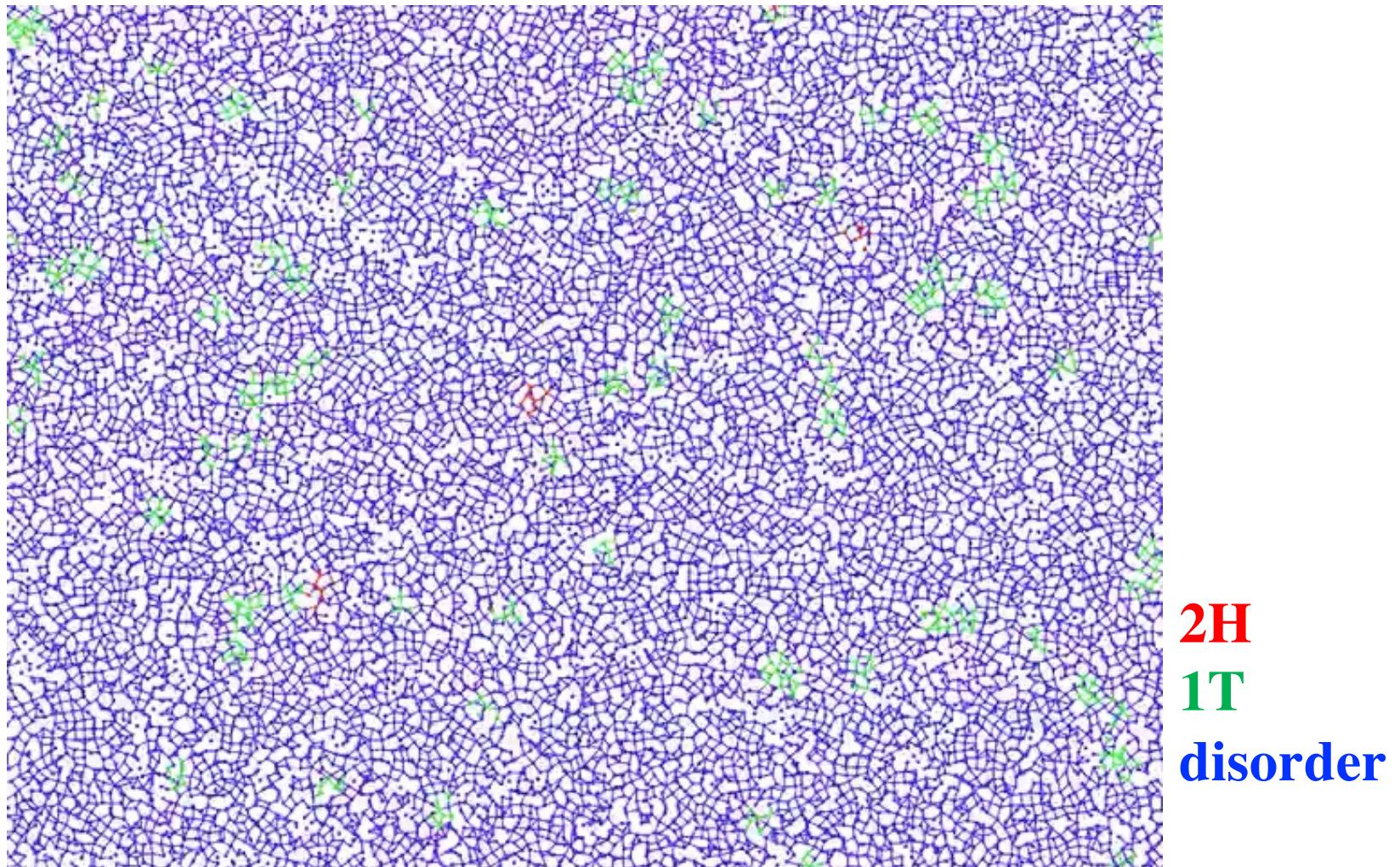
K. Liu *et al.*, Proc. ScalA18 ('18)  
S. Hong *et al.*, JPCl 10, 2739 ('19)

- Variational autoencoder to generate transformation pathways from images & latent-space algebra

P. Rajak *et al.*, Phys. Rev. B 100, 014108 ('19)



# Learning Transformation Pathways

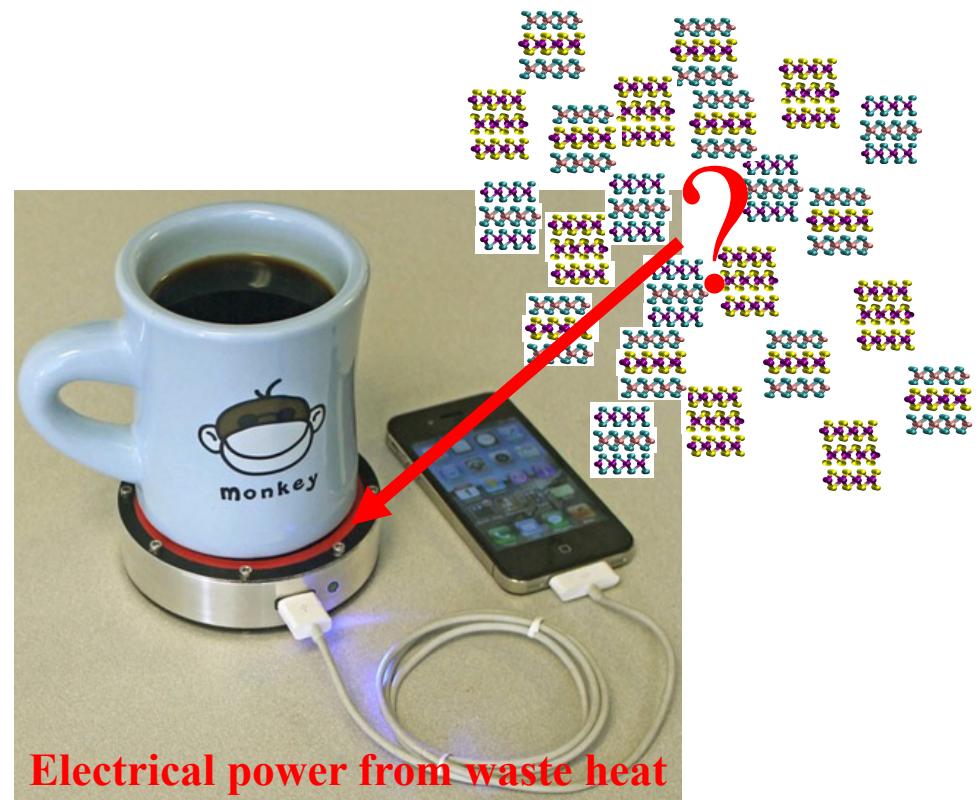
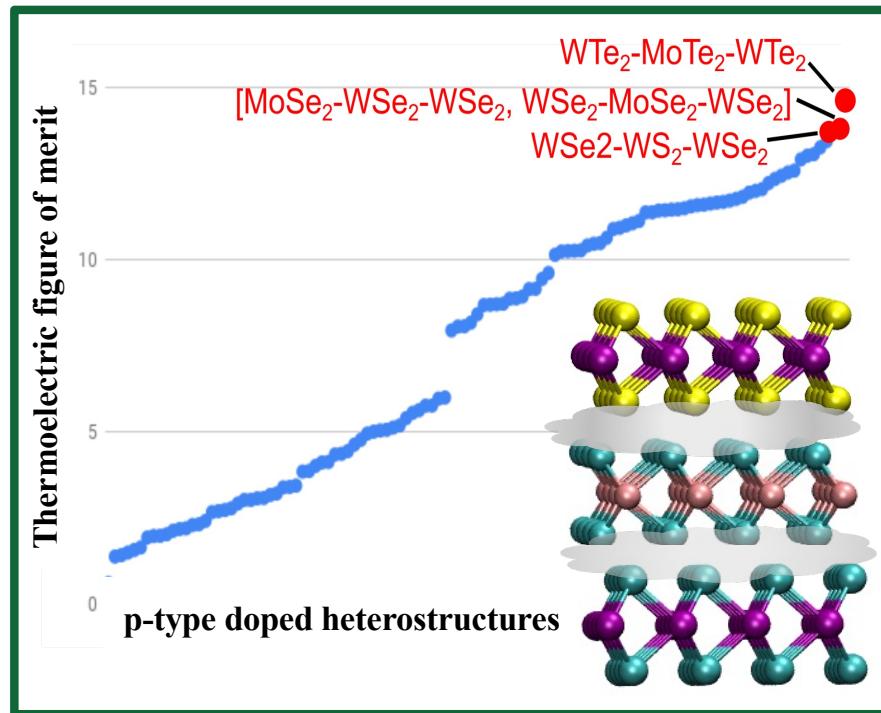


- Found novel transformation pathways to the stable 2H phase via the metastable 1T phase during chemical vapor deposition (CVD) growth of  $\text{MoS}_2$

S. Hong *et al.*, *J. Phys. Chem. Lett.* **10**, 2739 ('19)

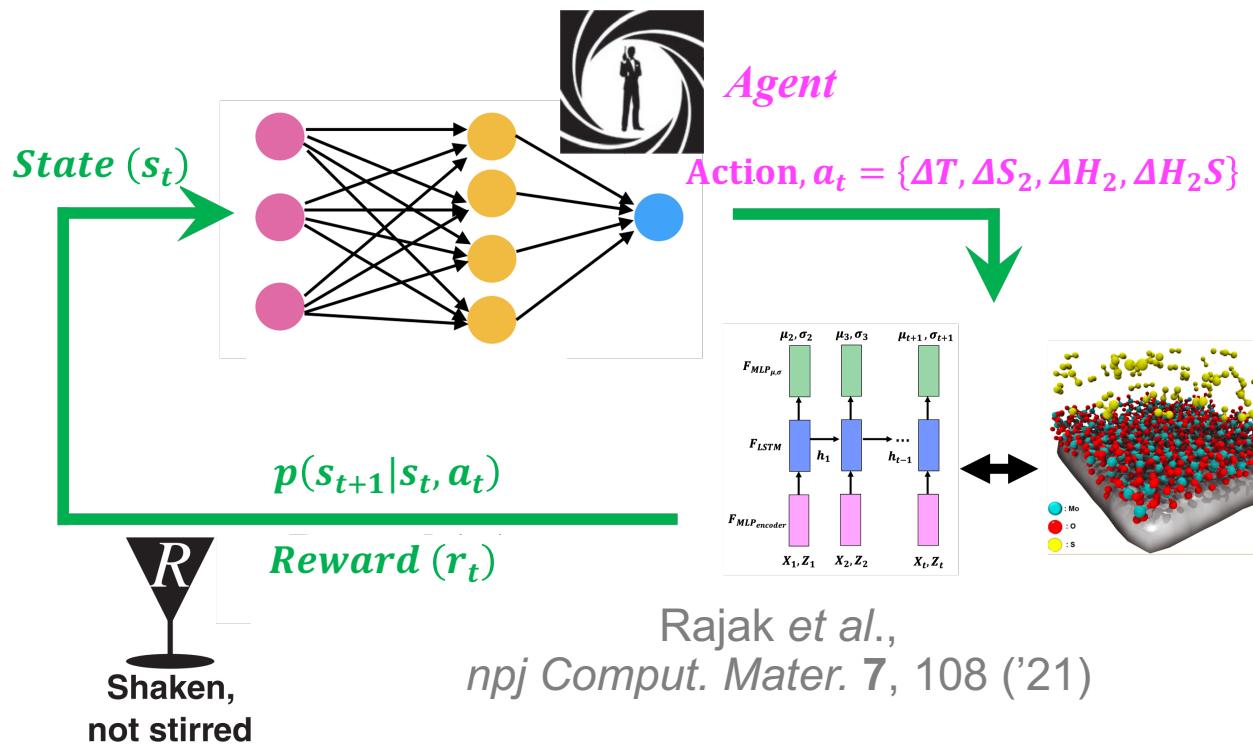
# Active Learning of Optimal Materials

- Bayesian optimization balances exploitation & exploration to find a structure with the desired property with a minimal number of quantum-mechanical calculations
- Predicted three-layered transition-metal chalcogenide (TMDC) heterostacks with the largest thermoelectric figure-of-merit



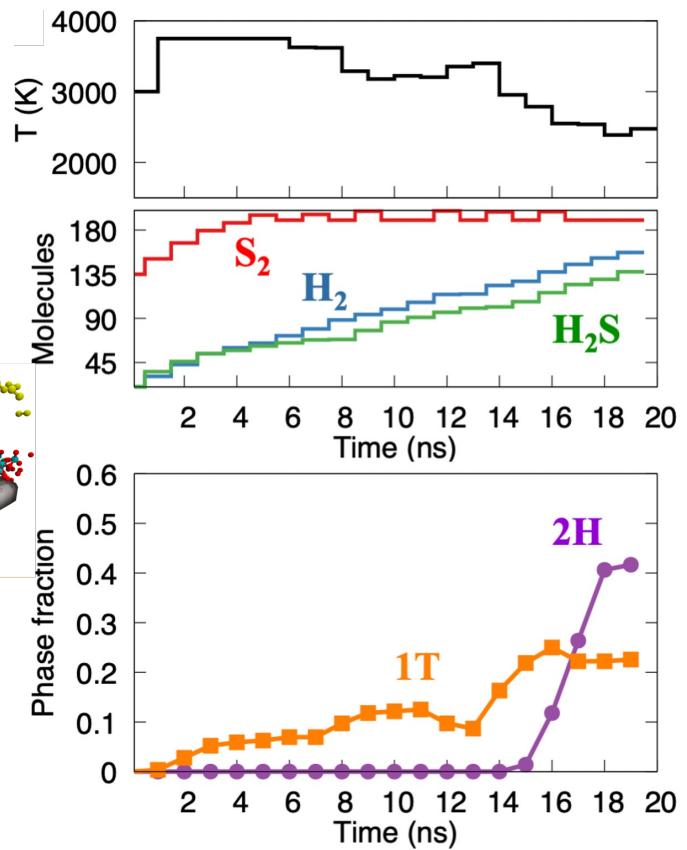
# Reinforcement Learning for Growth

- In a manner AI plays a board game of Go, use reinforcement learning (RL) to design optimal growth conditions (e.g., temperature & gas-pressure control) to achieve desired properties such as minimal defect density
- AI model combines:
  1. RL agent to design actions
  2. Neural network-based dynamic model trained by molecular-dynamics (MD) simulation to predict new states



Rajak et al.,  
npj Comput. Mater. 7, 108 ('21)

cf. Sgroi et al., Phys. Rev. Lett. 126, 020601 ('21)



# AI Meets Kirigami

- Reinforcement learning to design optimal kirigami with maximal stretchability

Rajak *et al.*, *npj Comput. Mater.* 7, 102 ('21)

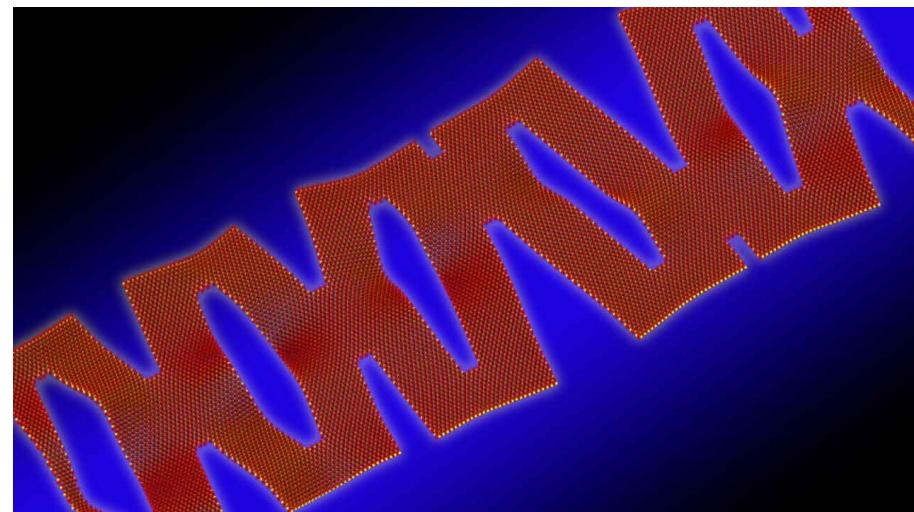
FEATURE STORY | ARGONNE NATIONAL LABORATORY

## Ancient art meets AI for better materials design

BY JOHN SPIZZIRRI | APRIL 7, 2022

Ancient Japanese art of kirigami guides artificial intelligence (AI) technique for durable, wearable electronics.

Kirigami is the Japanese art of paper cutting. Likely derived from the 剪纸 Chinese art of jiānzhǐ, it emerged around the 7<sup>th</sup> century in Japan,



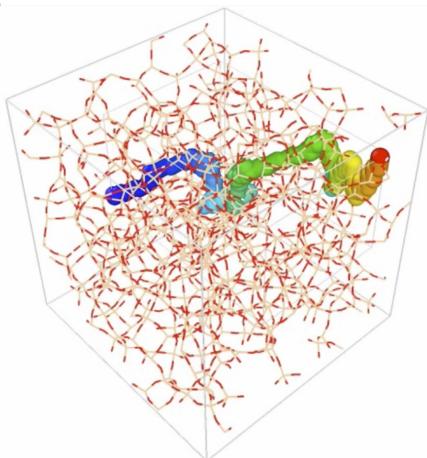
<https://www.anl.gov/article/ancient-art-meets-ai-for-better-materials-design>

# Reinforcement Learning for Long-Time Dynamics

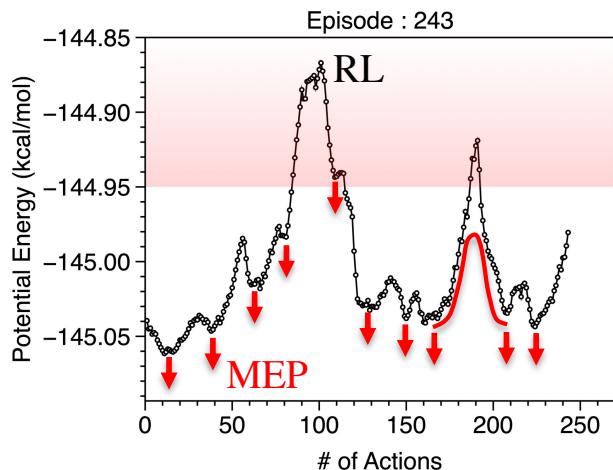
- Phase 1—explore (agent parallelism): Multiple reinforcement learning (RL) agents autonomously discover *long low-activation-barrier pathways*  
Mnih, *Nature* 518, 529 ('15); Hessel, *AAAI* 32, 11796 ('18)
- Phase 2—refine (time parallelism): Concurrent nudged-elastic-band (NEB) refinements of multiple *minimum-energy path (MEP) segments*
- Estimate time based on transition-state theory

$$t_{\text{migration}} = \sum_{i \in \{\text{activation events}\}} \frac{\hbar}{k_B T} \exp\left(\frac{E_i^{\text{activation}}}{k_B T}\right)$$

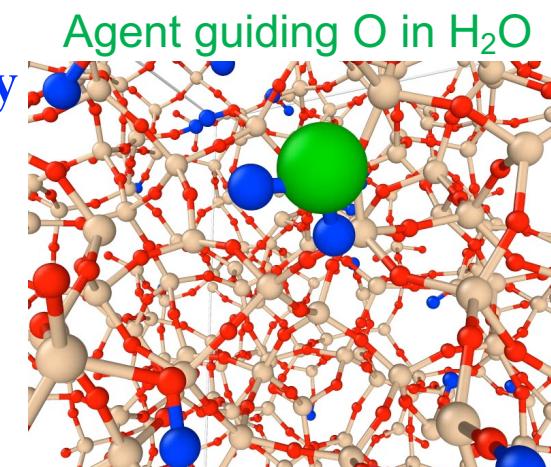
1. Explore



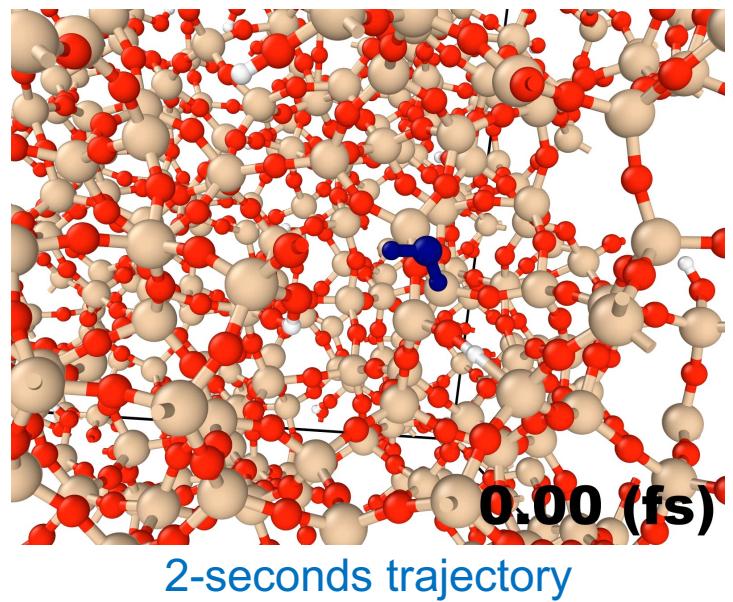
2. Refine



Nomura et al., *J. Phys. Chem. Lett.* 15, 5288 ('24)



H<sub>2</sub>O diffusion & reaction in SiO<sub>2</sub>



*AI for long time!*

# Dielectric Polymer Genome

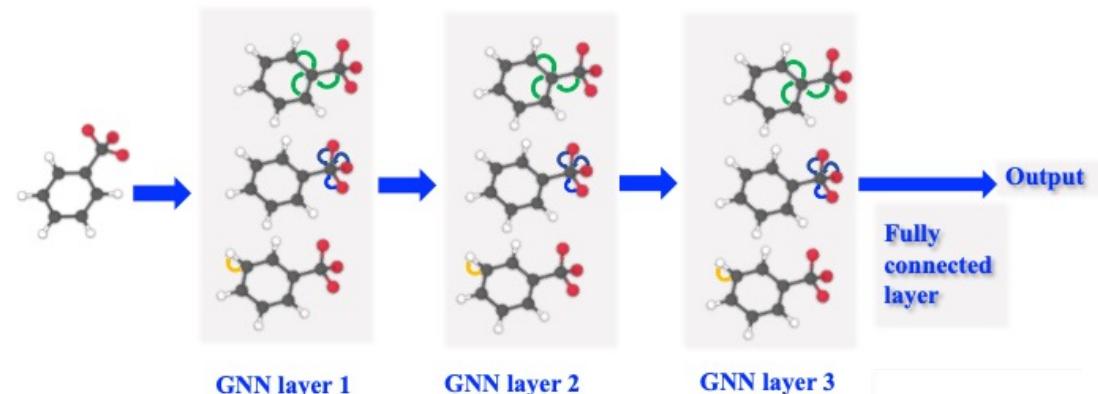
Recurrent neural network for polymer property prediction



May 2021 Volume 61, Issue 5

pubs.acs.org/jcim

Graph attention neural network for explainable property prediction



GNN layer 1

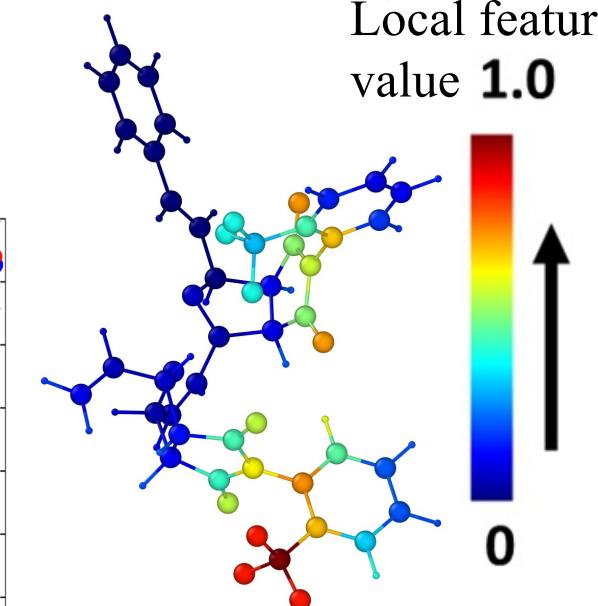
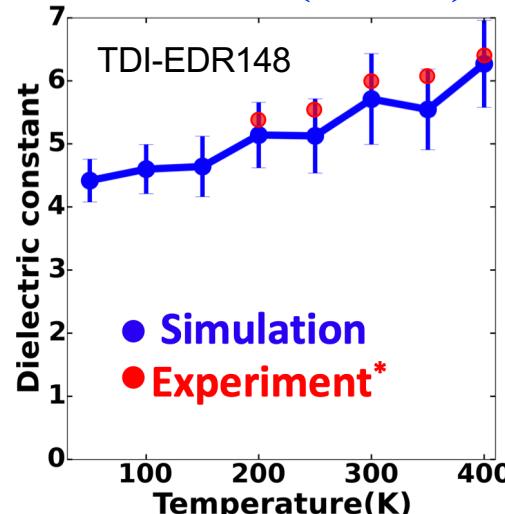
GNN layer 2

GNN layer 3

Fully  
connected  
layer

Local feature  
value **1.0**

Experimental  
validation (UConn)



Nazarova et al.,  
J. Chem. Info. Model. 61, 2175 ('21)

# Graph Neural Network

- **Allegro (fast) NNQMD: State-of-the-art *accuracy & speed* founded on *group-theoretical equivariance & local descriptors***

Musaelian et al., *Nat. Commun.* **14**, 579 ('23)

- An equivariant function  $f$  acts on an atomic geometry  $x \in X$  under a group action  $g$  (e.g., rotation) as  $f(g \cdot x) = g \cdot f(x)$

Thomas et al., *arXiv*. 1802.08219v3 ('18)

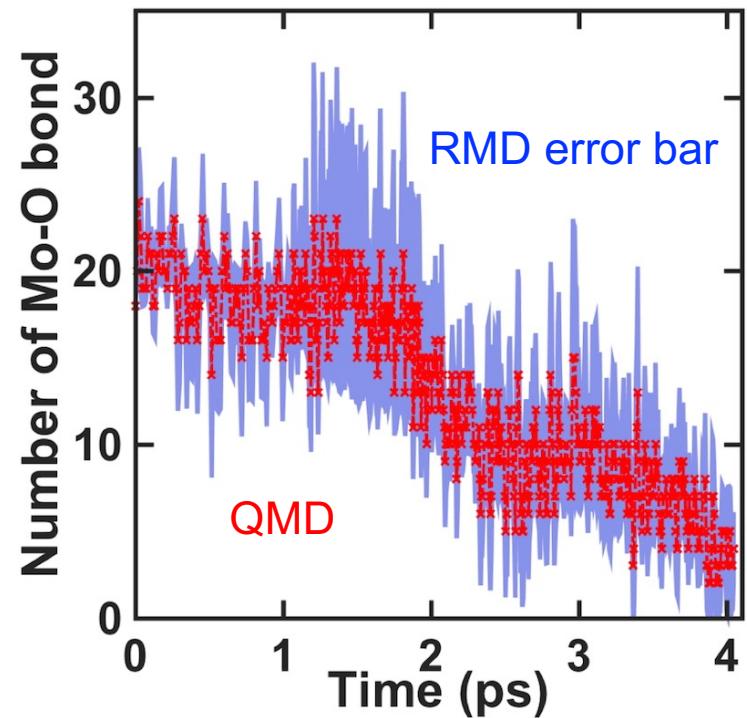
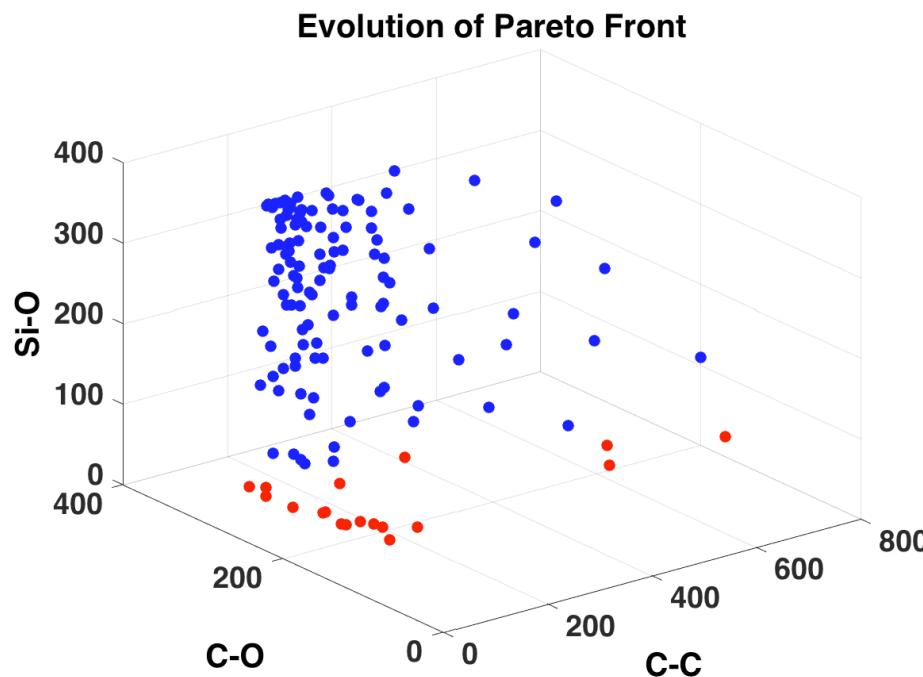
NNQMD: Neural network quantum molecular dynamics

$$\begin{array}{ccc} X & \xrightarrow{g \cdot} & X \\ f \downarrow & & \downarrow f \\ Y & \xrightarrow{g \cdot} & Y \end{array}$$

- Allegro speeds-up an equivariant message-passing (or graph neural-network) NNQMD, NequIP [Batzer et al., *Nat. Commun.* **13**, 2453 ('22)], by localizing atomic-environment descriptors, i.e., restricting message passing to the nearest neighbor
- A complete set of atomic-environment descriptors is provided by atomic cluster expansion (ACE) [Drautz, *Phys. Rev. B* **99**, 014104 ('19)]; relationship between NequIP & ACE has been discussed in Batatia et al., *arXiv*. 1802.08219v3 ('18)

# Pareto-Frontal Uncertainty Quantification

- Train reactive force-field parameters by dynamically fitting reactive molecular dynamics (RMD) trajectories to quantum molecular dynamics (QMD) trajectories on-the-fly
- Pareto optimal front in multiobjective genetic algorithm (MOGA) provides an ensemble of force fields to enable uncertainty quantification (UQ)



- Pareto-optimal solutions during genetic training (RMD errors for three quantities-of-interest)
- Converged Pareto-optimal front

A. Mishra *et al.*, *npj Comput. Mater.* **4**, 42 ('18)