



# Biomolecules Random Walks, Heterogeneities and Model Selection: *What Information is accessible from experimental Biomolecules Random Walk?*

Jean-Baptiste Masson

Physics of Biological Systems, Institut Pasteur, CNRS UMR 3525, Paris, France

Janelia Research Campus, Ashburn,

Laboratoire Physico-Chimie, Institut Curie, CNRS UMR168, Paris France

# Introduction

- Order of Magnitudes
- Analysing Random Walks
- Pipeline
- Not Tracking
- Mapping
- Application

# Biological Media

- Localization precision ~ 10nm
- Temporal precision ~ 10ms
- Static Density of tagging ~ 1-10000 #.μm<sup>-2</sup>
- Static Density of tagging ~ 1-100 #.μm<sup>-3</sup>
- Dynamic Density of Tagging ~ 0.01-10 #.μm<sup>-2</sup>
- Dynamic Density of Tagging ~ 0.01-2 #.μm<sup>-3</sup>
- Duration of recording ~ 1-1000s
- Field of View ~ 10-2500 μm<sup>2</sup>

- Diffusion ~ 0.005-5 μm<sup>2</sup>.s<sup>-1</sup>
- Drifts ~ 0.001-10 μm.s<sup>-1</sup>
- Forces ~ 0.001-10pN
- Interaction Energy ~ 0.1-10 k<sub>B</sub>T
- Local Gradients
  - $\nabla D \sim 10^{-3}$ -1 μm.s<sup>-1</sup>
  - $\nabla F \sim 1$ -500 k<sub>B</sub>T.μm<sup>-2</sup>
  - $\nabla V \sim 1$ -100 k<sub>B</sub>T.μm<sup>-1</sup>
  - $\nabla I \sim 0.001$ -10

## Time Scales:

- Transcription ~ 1 - 7.10<sup>5</sup> s
- Traduction ~ 1 - 7.10<sup>5</sup> s
- Methylation ~ 0.1 - 100 s
- Phosphorilation ~ 0.01 - 10 s
- Scaffold ~ 10 - 10<sup>5</sup> s
- Lipid Exchange ~ 0.1 - 500s
- Gephyrin ~ 300 s
- Un/binding ~ 0.1-100 s

## Spatial Scales:

- DNA ~ 0.3nm
- Actin-Myosin ~ 10nm
- Lipid Rafts ~ 20nm
- Caveolin ~ 50-70nm
- Lipid Platforms ~ 100-1000nm
- Scaffolds ~ 50-500nm
- PIP2/PIP3 ~ 200-3000nm
- Cytoskeleton Mesh ~ 50-1000nm
- Cytosol subdomain ~ 20-2000nm

# Analysing Random Walks

Physical Modelling  
Dictionary  
Features

Bayesian Inference  
Features ↔ Probability

Statistics  
Hypothesis  
Machine Learning  
Sampling

# Statistical Hypothesis

- Global distribution analysis (forced i.i.d.)
  - Heterogeneous
  - Statistical testing of pooling procedure, sub families?
- Single trajectory associated properties
  - Heterogeneous
  - Identical Environments?
- Clustered analysis
  - Spatial, internal states etc
  - Pooling neighbouring areas? Collapsing states?
- Time Evolution
  - Complex Windowing? Information?
  - Out of Equilibrium?

D. Koller and N. Friedman, Probabilistic Graphical Models: Principles and techniques

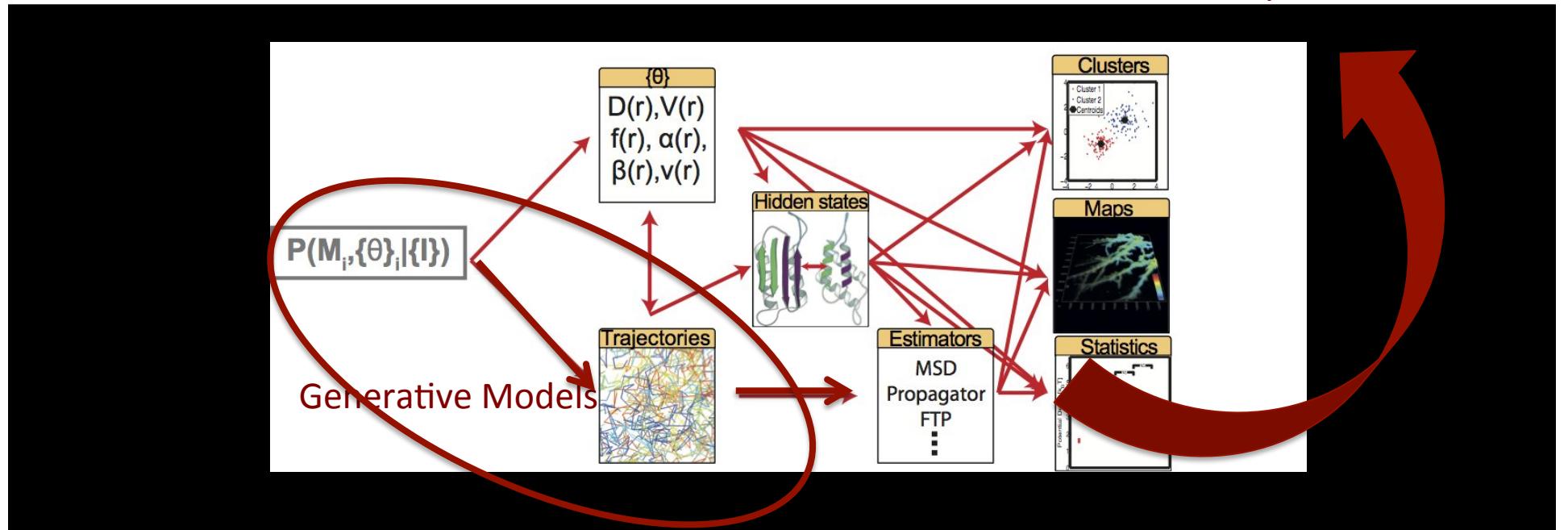
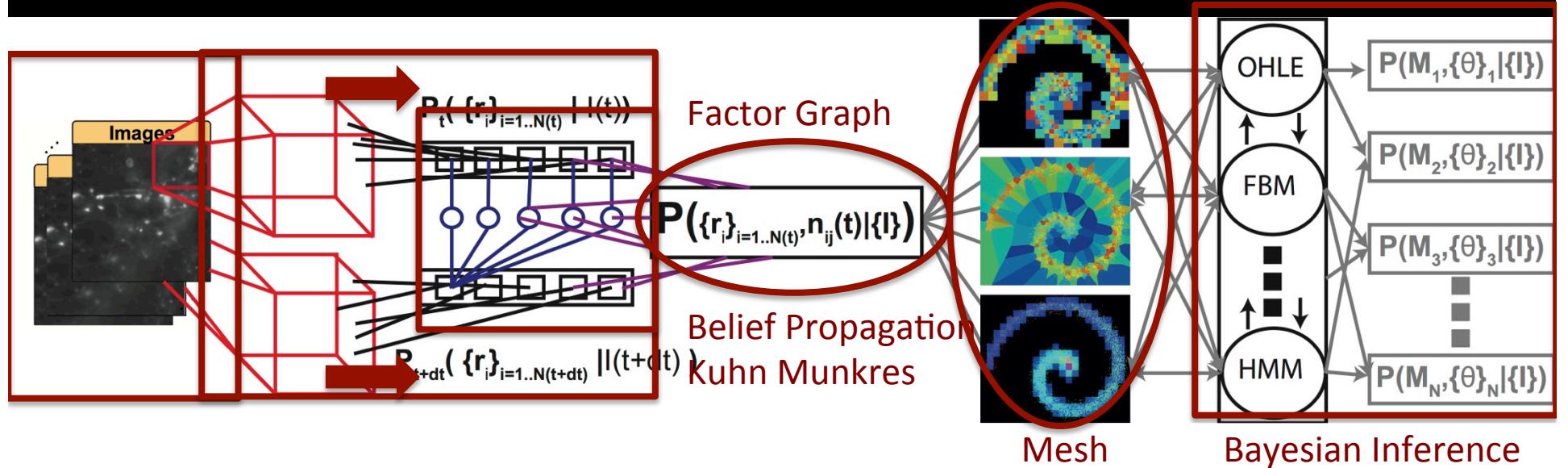
Y. Abu-Mostafa & M. Magdon-ismail, Learning from Data

C. Bishop, Pattern Recognition and Machine Learning

G. Saporta, Probabilités, Analyses de données et statistiques

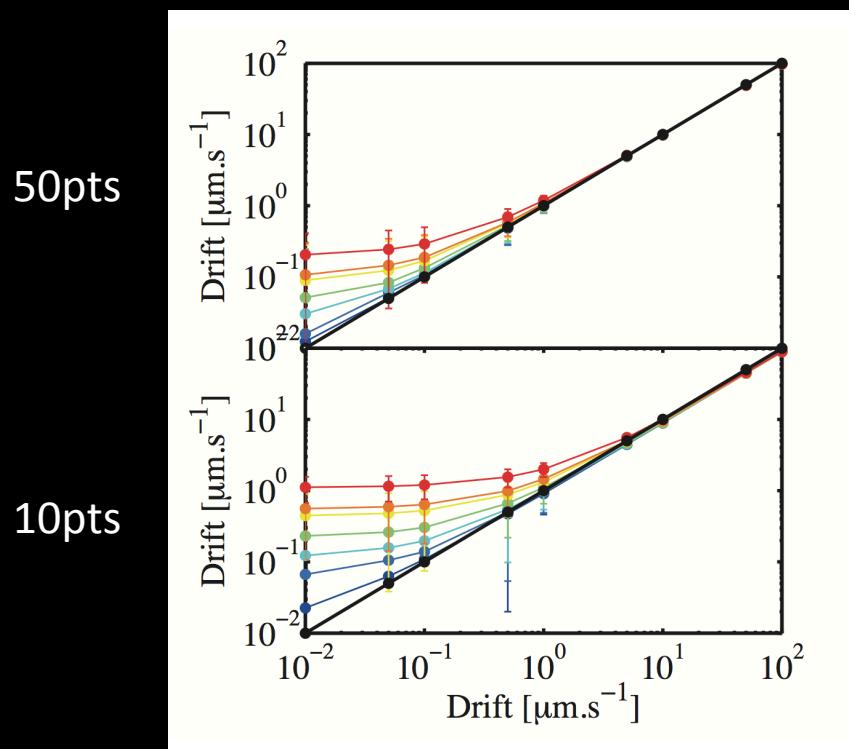
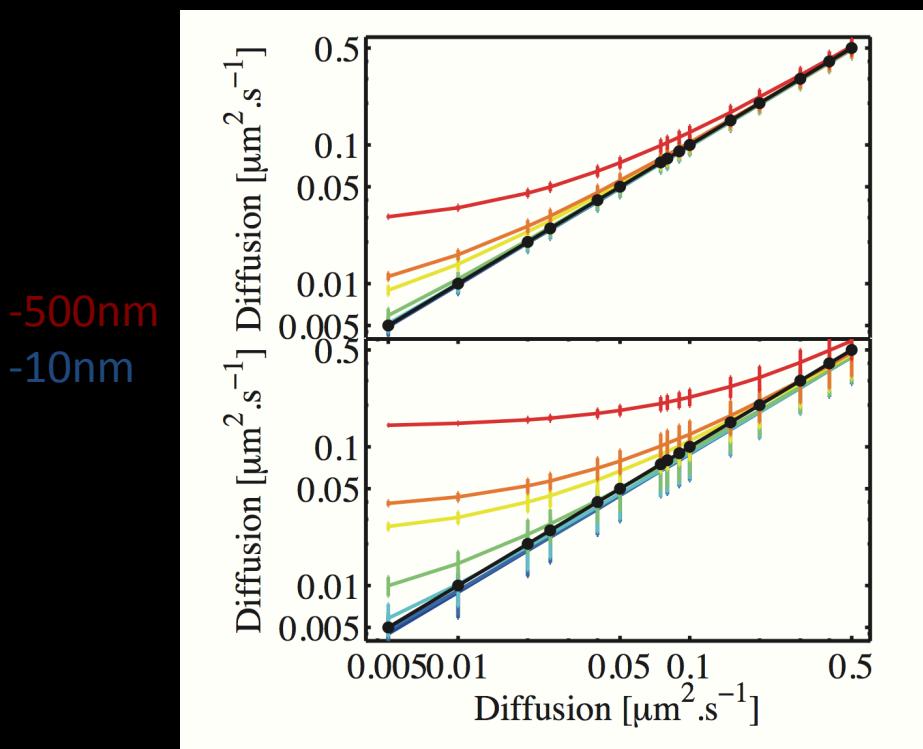
T. Hastie & R. Tibshirani , The Elements of Statistical Learning: Data Mining, Inference, and Prediction

# Pipeline



# Estimators and Acquisition Method

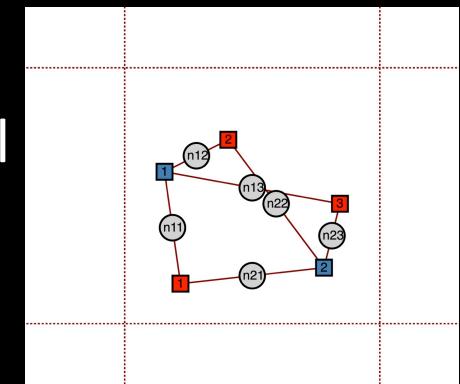
- Estimators have different level of robustness
- Acquisition methods matters
- Tracking error



- If Memoryless trajectories, High density trajectories

# No Tracking-Graph Matching: 2 Choices

- High Density ( $\rho < 5 \text{ #.} \mu\text{m}^{-2}$ ,  $\rho < 2 \text{ #.} \mu\text{m}^{-3}$ ) or highly variable number of Biomolecules
  - Best Assignment captures most of the dynamics
  - 2 steps
    - Min-Sum algorithm on distance based model
    - Bayesian Inference of the Parameters
- Very High Density ( $\rho > 5 \text{ #.} \mu\text{m}^{-2}$ ,  $\rho > 2 \text{ #.} \mu\text{m}^{-3}$ )
  - Tracking should not be performed
  - 1 step
    - Solve simultaneously the assignment and the inference Problem: Belief Propagation



# High Density Step 1: Min-Sum

- Distance based model

$$P(\{n\}|\theta) = \prod_{i=1}^{N_A} \left[ \delta\left(\sum_{j=1}^{N_B} n_{ij}, 1\right) + e^{-\mu} \delta\left(\sum_{j=1}^{N_B} n_{ij}, 0\right) \right] \times \prod_{j=1}^{N_B} \left[ \delta\left(\sum_{i=1}^{N_A} n_{ij}, 1\right) + e^{-\mu} \delta\left(\sum_{i=1}^{N_A} n_{ij}, 0\right) \right] \times \prod_{(i,j)} e^{-\beta n_{i,j} d_{ij}^2}$$

- $\beta \rightarrow \infty$ , only the best Assignment remains
- Iteratively find the best Assignment using Min-Sum Scheme

$$\begin{cases} x_{k \rightarrow i}^j = \min_{l \in I(j) \setminus i} \left( -\log(d_{l \rightarrow k}^{j \rightarrow j+1}) - x_{l \rightarrow k}^{j+1} \right) \\ x_{i \rightarrow k}^{j+1} = \min_{l \in I(j+1) \setminus k} \left( -\log(d_{i \rightarrow l}^{j \rightarrow j+1}) - x_{l \rightarrow i}^j \right) \end{cases} \quad \begin{cases} \pi^j(i) = \operatorname{argmin}_l \left( -x_{l \rightarrow k}^{j+1} - \log(d_{i \rightarrow l}^{j \rightarrow j+1}) \right) \\ \pi^{j+1}(k) = \operatorname{argmin}_i \left( -x_{i \rightarrow k}^j - \log(d_{i \rightarrow k}^{j \rightarrow j+1}) \right) \end{cases}$$

- With

$$\pi^j(\pi^{j+1}(i)) = i$$

M. Mezard & A. Montanari, Information, Physics , and computation  
Yedidia et al, IEEE TRANS. on inf. theo., vol. 51, no. 7, (2005)

# Very High density #10. $\mu\text{m}^{-2}$

$$P(\{n\}|\theta) = \prod_{i=1}^{N_A} \left[ \delta\left(\sum_{j=1}^{N_B} n_{ij}, 1\right) + e^{-\mu} \delta\left(\sum_{j=1}^{N_B} n_{ij}, 0\right) \right] \times \prod_{j=1}^{N_B} \left[ \delta\left(\sum_{i=1}^{N_A} n_{ij}, 1\right) + e^{-\mu} \delta\left(\sum_{i=1}^{N_A} n_{ij}, 0\right) \right] \times \prod_{(i,j)} e^{-\beta n_{i,j} \sqrt{-4\pi(D+\frac{\sigma^2}{\delta t})\delta t - \frac{d_{ij}^2}{4(D+\frac{\sigma^2}{\delta t})\delta t}}}$$

$$\{(D_l)\} = \underset{(D_l)}{\operatorname{Arg\,max}} \sum_{\{n\}} P(\{n\}|(D_l)) = \underset{(D_l)}{\operatorname{Arg\,max}} Z(D_l) = \underset{(D_l)}{\operatorname{Arg\,max}} \log(Z(D_l)) = \underset{(D_l)}{\operatorname{Arg\,min}}(F)$$

- Approximating the Free Energy: Bethe Free Energy

$$b(\{n\}) \propto \frac{\prod_i b_i(n_i) \prod_j b_j(n_j)}{\prod_{(i,j)} b_{ij}(n_{ij})}$$

$$x_{j \rightarrow i}^L = -\frac{1}{\beta} \log \left( \sum_{k \in A \setminus i} e^{\log(P_{kj}) + \beta x_{k \rightarrow j}^R} + e^{-\mu} \right)$$

$$x_{i \rightarrow j}^R = -\frac{1}{\beta} \log \left( \sum_{k \in B \setminus j} e^{\log(P_{ik}) + \beta x_{k \rightarrow i}^L} + e^{-\mu} \right)$$

$$F(\beta) = -\sum_{(i,j)} \log \left( 1 + e^{\log(P_{kj}) + \beta x_{i \rightarrow j}^R + \beta x_{j \rightarrow i}^L} \right)$$

$$+ \sum_{i \in A} \log \left( e^{-\mu} + \sum_{j \in B} e^{\log(P_{kj}) + \beta x_{j \rightarrow i}^L} \right)$$

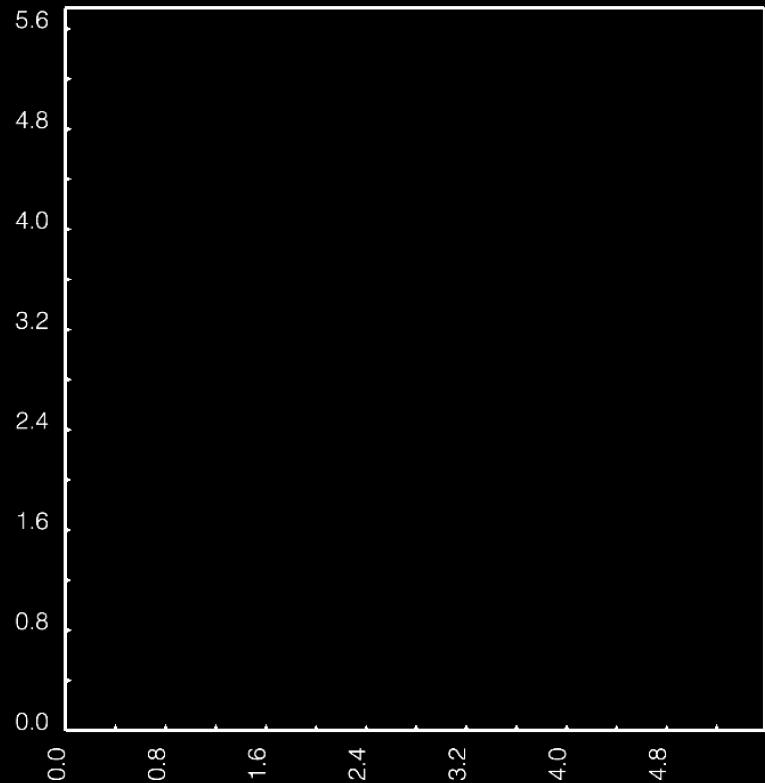
$$+ \sum_{j \in B} \log \left( e^{-\mu} + \sum_{i \in A} e^{\log(P_{kj}) + \beta x_{i \rightarrow j}^R} \right)$$

# Selective Graph Cuts

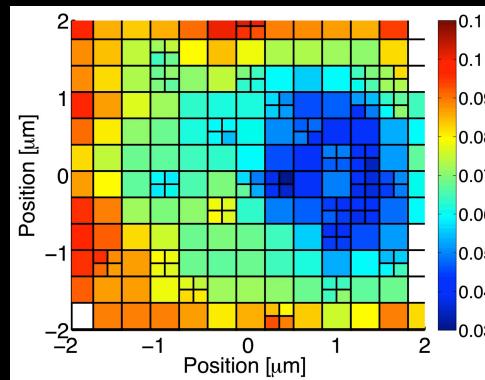
- Highly spatially varying parameters
- Global Solution overestimate parameters and induce high spatial smoothing
- Recursive Solving
  - Until Convergence (or predefined )
  - Parameters initialisation (Over estimation)
  - Mesh generation (Hierarchical)
  - Graph Cuts Based on the Mesh (  $O(N)$  )
  - Belief Propagation to extract parameters

# Temporal evolution of diffusive trap

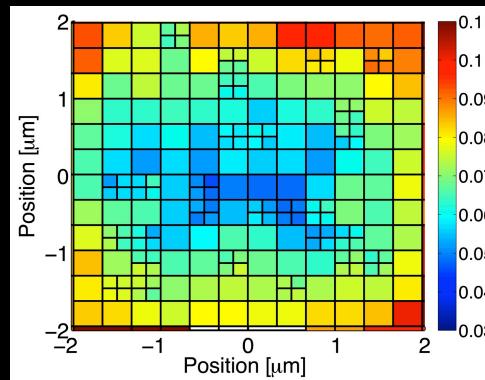
**Double Trap**



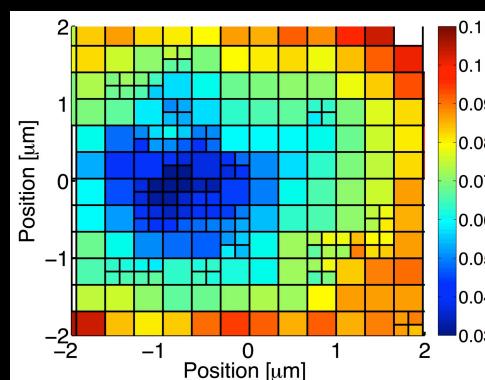
Numerical Simulations  
Scaffolding proteins evolution



$t = 0\text{ s}$



$t = 30\text{ s}$



$t = 60\text{ s}$

# Bayesian Inference

$$P(\{U_k\} \mid \{T_l\}, M_i) = \frac{P(\{T_l\} \mid \{U_k\}, M_i) P(\{U_k\} \mid M_i)}{P(\{T_l\} \mid M_i)}$$

- $P(\{U_k\} \mid \{T_l\}, M_i)$  posterior distribution of the parameters  $\{U_k\}$  having observed  $\{T_l\}$  (set of trajectories) with Model  $M_i$
- $P(\{T_l\} \mid \{U_k\}, M_i)$  likelihood of model
- $P(\{U_k\} \mid M_i)$  prior distribution of the parameters
- $P(\{T_l\} \mid M_i)$  Evidence of the Model
- Bayesian Choice 
$$\gamma_{ij} = \frac{P(\{T_l\}, M_i)}{P(\{T_l\}, M_j)} = \frac{P(\{T_l\} \mid M_i) P(M_i)}{P(\{T_l\} \mid M_j) P(M_j)}$$
- Bayesian Choice Systematically fails for most Random Walks analysis

# Biomolecules Random Walk

## Overdamped Langevin Equation

Interaction Terms

Diffusive Term

$$\frac{d\vec{r}}{dt} = \frac{F_t(\vec{r})}{\gamma_t(\vec{r})} + \sqrt{2D_t(\vec{r})}\xi(\vec{t})$$

$$F_t(\vec{r}) = -\nabla V(\vec{r})$$

$$\gamma_t(\vec{r}) = \frac{k_B T}{D_t(\vec{r})}$$

$$\langle \xi_i(t)\xi_j(s) \rangle = \delta(t-s)\delta_{ij}$$

- The Fokker-Planck Equation, Methods of solutions and Applications, H. Risken
- J.-B M et al, Phys. Rev. Lett 102 , 048103 (2009)
- J.-B M et al, Biophys J. vol 106, issue 1, p74-83 (2014)
- Hanggi et al, J. Stat. Phys. Vol 18, No. 2 (1978)
- O. Farago & N. Grønbech-Jensen, J. Stat. Phys. Vol 156, p1093 (2014)

- O. Farago & N. Grønbech-Jensen, Phys. Rev. E 89, 013301 (2014)
- R. Fox, J. Stat. Phys. Vol 46 issue 5/6 (1987)

# Priors

- Experimental Priors: statistics from previous experiments
- Simulation Priors: Statistics and Cost Functions

C. Salvatico *et al*, Nat. Neuro. (submitted)

- Jefferey's Prior

JBM *et al*, Biophys. J.,  
vol. 106 issue 1 (2014)

$$P(\{U\}) \propto \sqrt{|J|} \quad J = \left( \partial_U \partial_{U'}^T \int d\mathbf{r} \sqrt{P(\mathbf{r}| \{U\}) P(\mathbf{r}| \{U'\})} \right)_{U=U'}$$
$$P(D(\vec{r})) \propto \frac{D^2(\vec{r})}{(D(\vec{r}) \Delta t + \sigma^2)^2}$$

- Field Priors

$$P(D(\vec{r}), V(\vec{r})) \propto \exp \left( -\mu \iint \|\nabla D(\vec{r})\|^2 d^2 r - \lambda \iint \|\nabla V(\vec{r})\|^2 d^2 r \right)$$

- Non-Local Priors  
optimization)

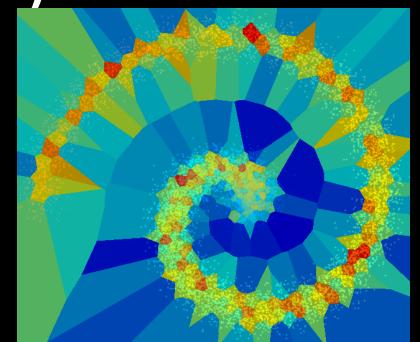
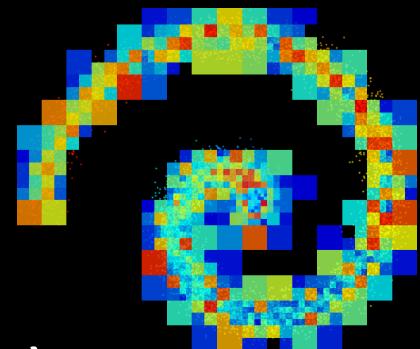
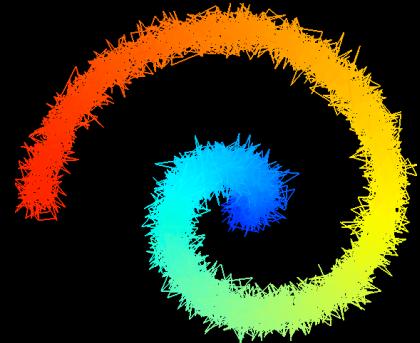
$$P(D(\vec{r})) \propto \exp(-\sigma(D(\vec{r}), \vec{F}(\vec{r}))) \quad (\text{complex})$$

# Multi-Scale Heterogeneities

- Unknown level or scale of heterogeneities
- Level :  $\Delta D \sim 0.01 - 0.2 \text{ } \mu\text{m}^2.\text{s}^{-1}$  (100nm) synapses
- Scales :  $l \sim 1\text{nm}-1\mu\text{m}$ . Full Cell Scale.
- Measures at Small Scales:  $dt \sim 10\text{ms}$  ,  $dr \sim 10\text{nm}$
- Partially Unsupervised Learning Problem
- Local single Molecule Density might be a clue (Biology)

# 2 Main Methods:

- Quadtree (Octree for 3D)
  - Hierarchical Graph
  - Subdivide space
  - 2 criteria
    - Capacity -> Information
    - Length scale -> local diffusivity scale
- K-Mean Clustering (Gaussian Mixture,...)
  - Unsupervised Learning
  - Initialization (Random, Density Based, Bubbling)
  - K-Mean clustering
  - Hierarchical updates (divide, fusion etc)



# Fokker-Planck and the Ito/Stratonovitch/Hanggi Dilemma

- Different Definition of the Stochastic Integrals

$$\frac{\partial P(\mathbf{r}, t | \mathbf{r}_0, t_0)}{\partial t} = -\nabla \cdot \left[ \left( \frac{F_t(\mathbf{r})}{\gamma_t(\mathbf{r})} + \lambda \nabla D_t(\mathbf{r}) \right) P(\mathbf{r}, t | \mathbf{r}_0, t_0) \right] + \nabla \cdot [D_t(\mathbf{r}) \nabla P(\mathbf{r}, t | \mathbf{r}_0, t_0)]$$

- Ito  $\lambda=0$ , Stratonovitch  $\lambda=1/2$  and Hanggi  $\lambda=1$
- $D(\mathbf{r} \rightarrow \mathbf{r}_0) \sim D(\mathbf{r}_0)$ ,
- $F(\mathbf{r} \rightarrow \mathbf{r}_0) \sim F(\mathbf{r}_0)$ ,
- $\nabla V(\mathbf{r} \rightarrow \mathbf{r}_0) \sim \nabla V(\mathbf{r}_0)$
- $D(\mathbf{r} \rightarrow \mathbf{r}_0) \sim D(\mathbf{r}_0) + \nabla D \cdot (\mathbf{r} - \mathbf{r}_0)$
- $F(\mathbf{r} \rightarrow \mathbf{r}_0) \sim F(\mathbf{r}_0)$ ,
- $\nabla V(\mathbf{r} \rightarrow \mathbf{r}_0) \sim \nabla V(\mathbf{r}_0)$

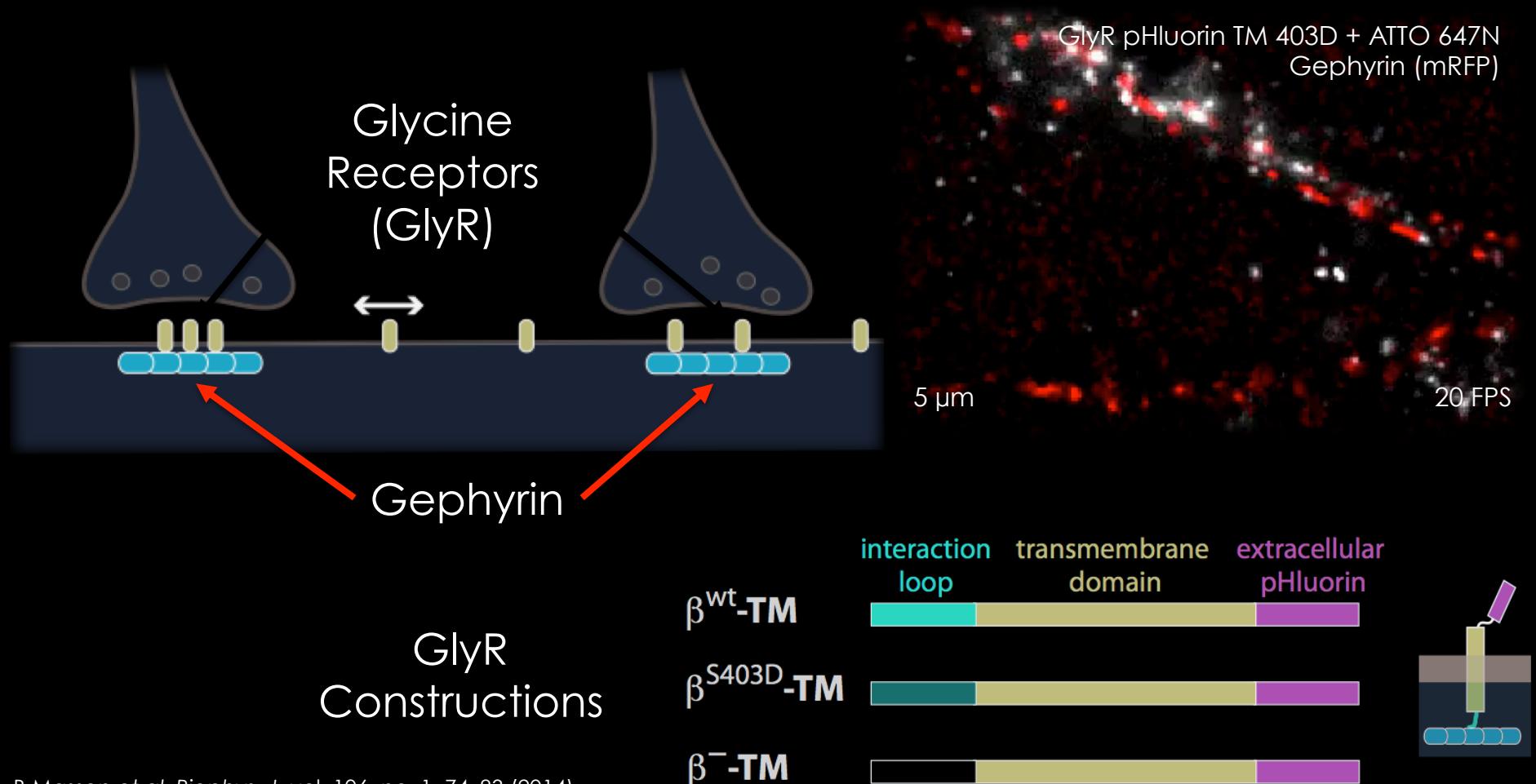
$$P(\mathbf{r}, t; \mathbf{r}_0, t_0 | D, F) = \frac{\exp\left(-(\mathbf{r} - \mathbf{r}_0 - DF\Delta t/k_B T)^2 / 4\left(D + \frac{\sigma^2}{\Delta t}\right)\Delta t\right)}{4\pi\left(D + \frac{\sigma^2}{\Delta t}\right)\Delta t}$$

$$P(\mathbf{r}, t; \mathbf{r}_0, t_0 | D, \nabla D, F) = \frac{1}{2\pi} \int_{-\infty}^{\infty} \frac{1}{(1 + ik\nabla D\Delta t)^{\left(\frac{F/\gamma + (1+\lambda)\nabla D}{\nabla D}\right)}} e^{-\frac{k^2\Delta t\left(D + \sigma^2/\Delta t - \nabla D x_0\right) + ikx_0}{1 + ik\nabla D\Delta t}} e^{ikx}$$

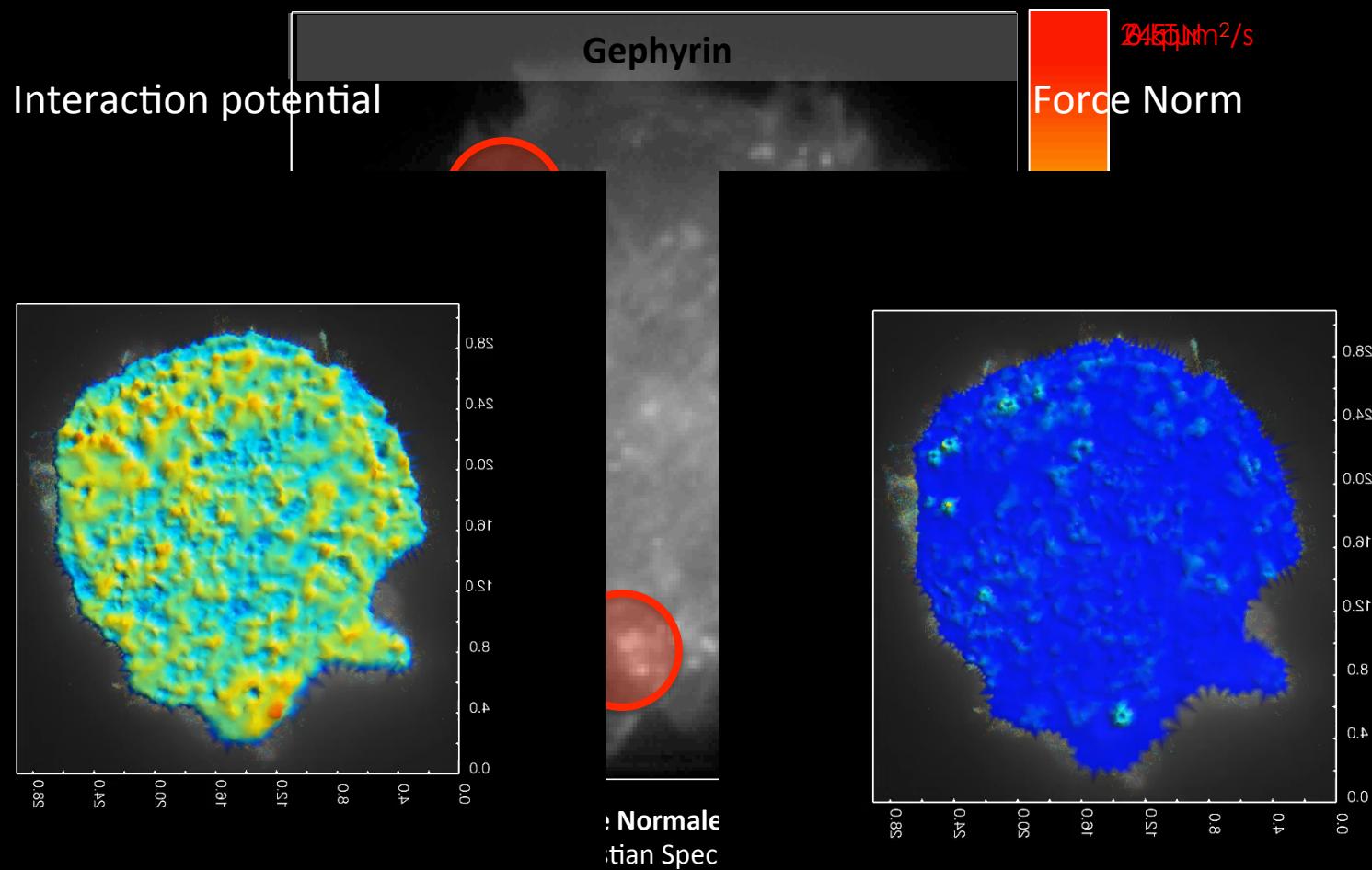
The Fokker-Planck Equation, Methods of solutions and Applications, H. Risken  
 O. Farago, N. Grønbech-Jensen, J. Stat. Phys. 156, 1093, (2014)  
 P. Hanggi, Helvetica Physica Acta, Vol. 51 (1978)  
 R. L. Fox, J. Stat. Phys., Vol. 46, Nos. 5/6, 1987

# Inhibitory synapses: GlyR and Gephyrin

- Gephyrin is a critical scaffold protein in the inhibitory neuronal synapse
- Synapses in neurons and artificial



# Artificial Synapses in Cos Cells



# Generative Model: Simulation in The Maps

- Fokker Planck  $\longrightarrow$  Master Equation
- Masters Equations

$$\frac{dP_{(i,j)}(t)}{dt} = \sum_{(i',j') \in N(i,j)} W_{(i,j),(i',j')} P_{(i',j')}(t) - \sum_{(i',j') \in N(i,j)} W_{(i',j'),(i,j)} P_{(i,j)}(t) \quad a_v = W_{(i,j),(i',j')}$$
$$W_{(i,j),(i',j')} = \frac{D(i',j')}{\Delta x^2} \exp\left(-\frac{\Delta x F_{(i,j),(i',j')}^x}{2\gamma(i',j') D_{(i',j')}}\right) \quad a_0 = \sum_v a_v$$

- Multi-scale Space Structure

- Time to wait

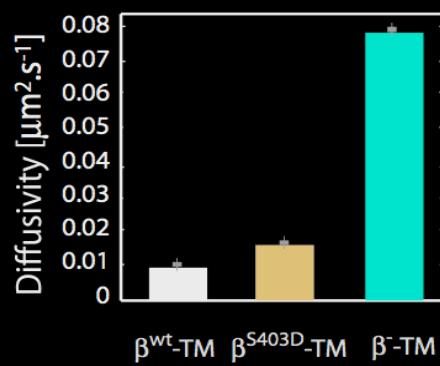
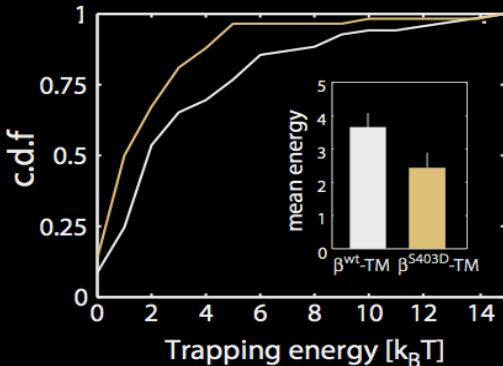
$$\tau = \frac{1}{a_0} \log\left(\frac{1}{r_1}\right) \quad r_1 = [0..1]$$

- Site to chose, k,

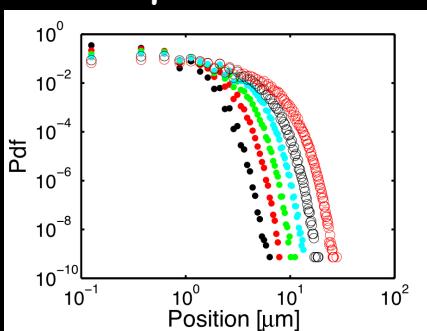
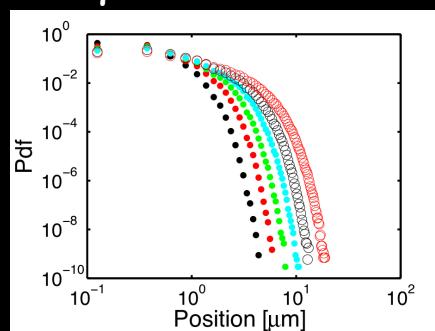
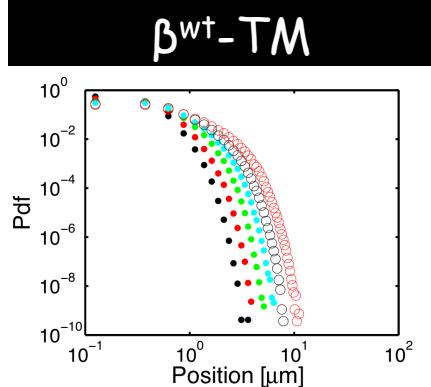
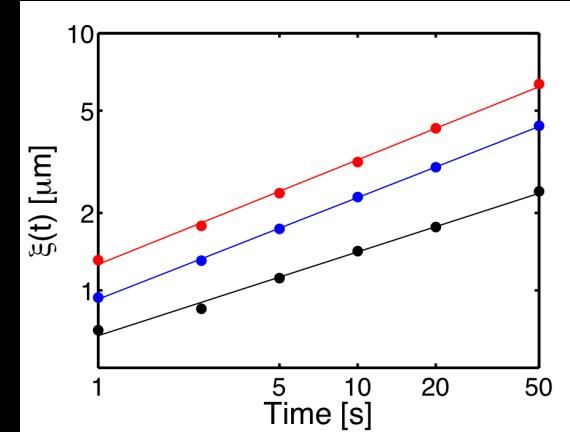
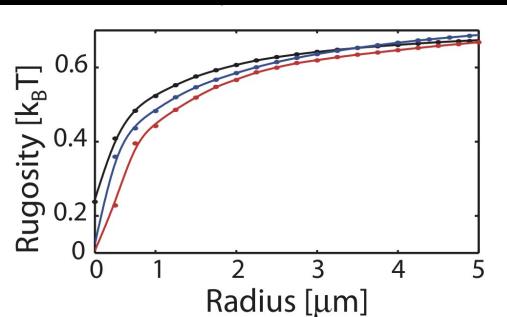
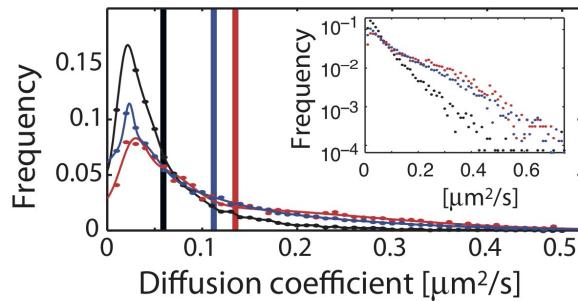
$$\sum_{v=0}^{k-1} a_v \leq r_2 a_0 < \sum_{v=0}^k a_v \quad r_2 = [0..1]$$

# Few Results

In the synapses



$\beta^{WT-TM}$ ,  $\beta^{S403D-TM}$ ,  $\beta^--TM$



$$\xi(t) \propto t^\gamma$$

$\beta^--TM$ :  $\gamma=0.41$

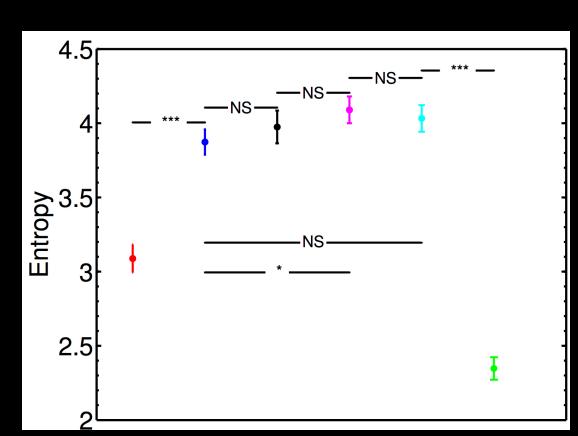
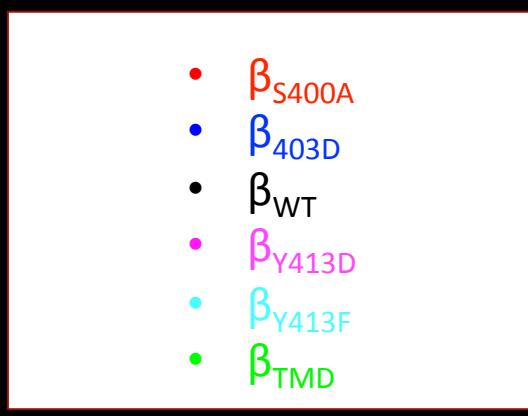
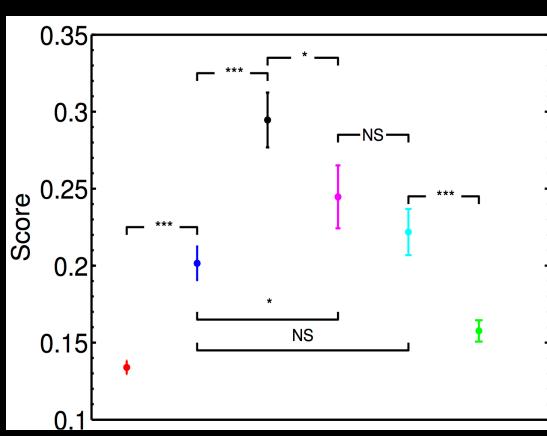
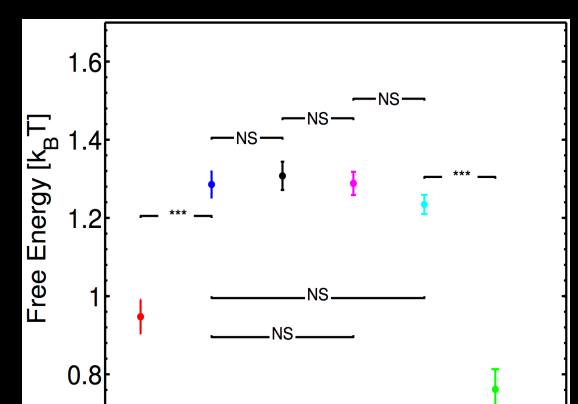
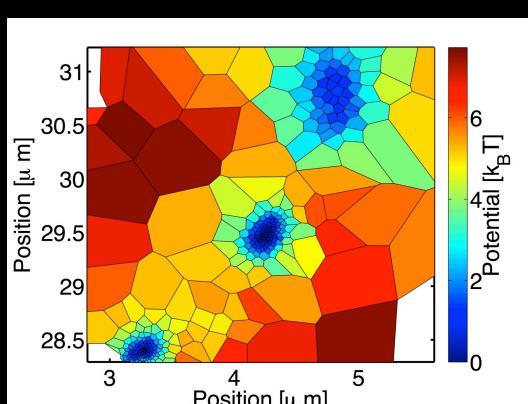
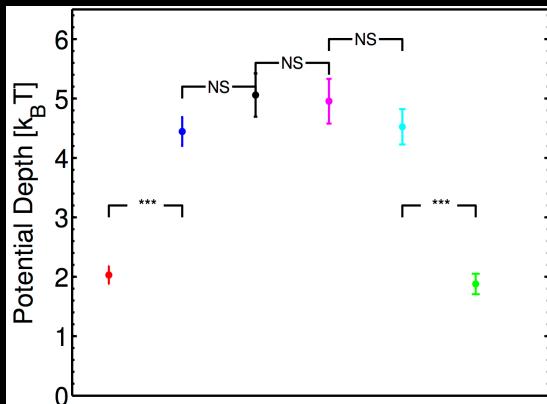
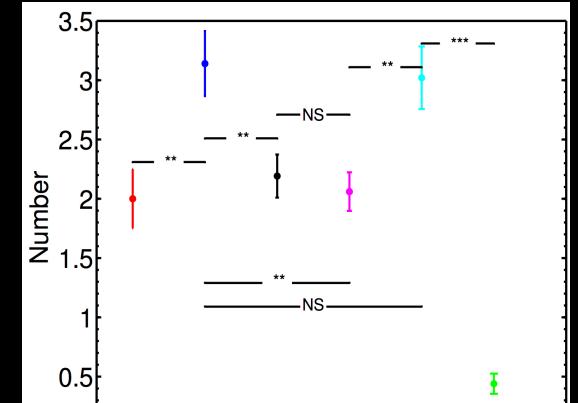
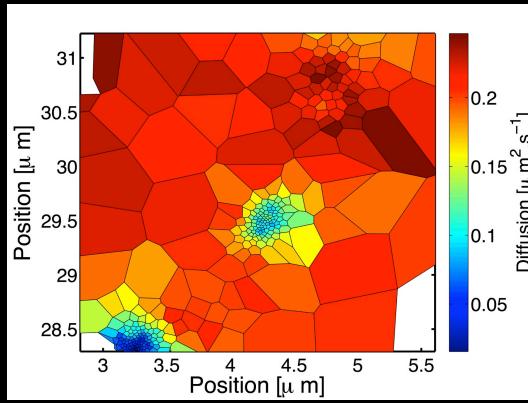
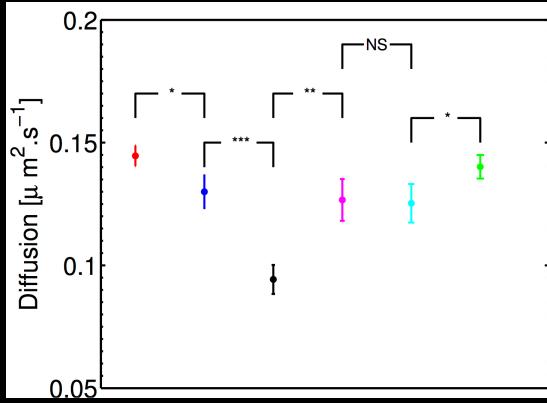
$\beta^{S403D-TM}$ :  $\gamma=0.39$

$\beta^{WT-TM}$ :  $\gamma=0.33$

$$P(r) = \frac{1}{\pi \xi^2(t)} e^{-\frac{r^2}{\xi^2(t)}}$$

$$\rho = \frac{r}{\xi(t)}$$

# Pipeline On Neurons: 450 Go $\rightarrow$ 670 synapses (1 Hour)

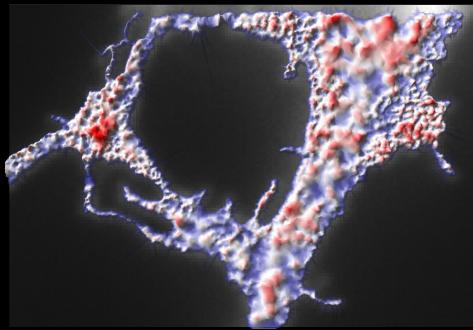


# Equilibrium?

- Is the process at equilibrium?
- Methods to detect Characteristics of Non-Equilibrium dynamics
  - Work of the force fields  $\Delta W = \oint \vec{F} \cdot d\vec{r}$  on random close contour in the Maps  $\vec{F} = -\vec{\nabla}V + \vec{\nabla} \times \vec{A}$  Ill posed problem
  - Asymmetries in the paths (A->B, B->A)
    - Metric on the paths
    - Relative Entropy
  - Detailed Equilibrium measure at various scales of the mesh  $\chi = P_i P_{i \rightarrow j} / P_j P_{j \rightarrow i}$
- All compared to simulations (information accumulated)  
R. Kawai *et al.*, PRL 98, 080602 (2007) M. Esposito *et al.*, N. J. Phys., 12 013013 (2010)  
A. Comez-Marin *et al.*, EPL 82 50002 (2008) J. Horowitz & C. Jarzynski, PRE 79 021106 (2009) A. V. Popov & R. Hernandez, PRE

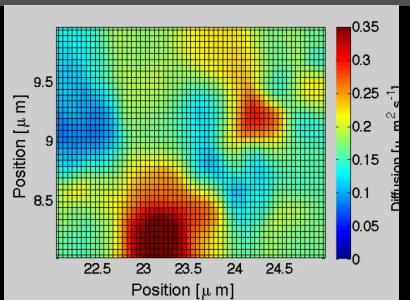
# InferenceMAP

Transmembrane Protein Diffusion  
in Neurons



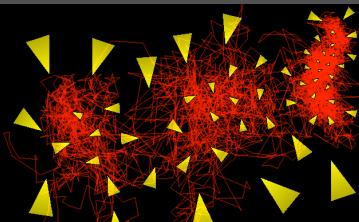
P. Dionne, Ecole Normale Supérieure

Rac<sub>1</sub> dynamics



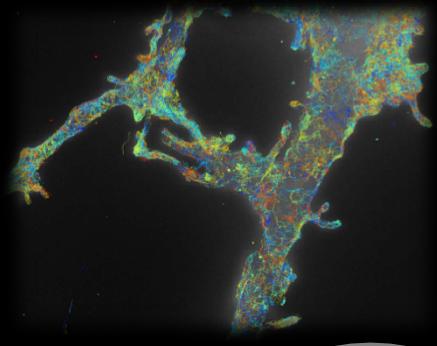
A. Remorino, Institut Curie

Membrane Confinement Zones

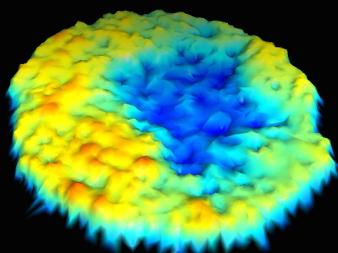


S. Türkcan, Ecole Polytechnique

Tailored for Big Data



Nuclear Protein Diffusion



Calculation Flexibility

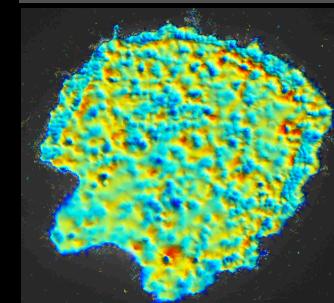
$$P(D(\vec{r}), \vec{F}(\vec{r}))|T)$$

$$P(D(\vec{r}), V(\vec{r}))|T)$$

J. Liu, Janelia Research Campus

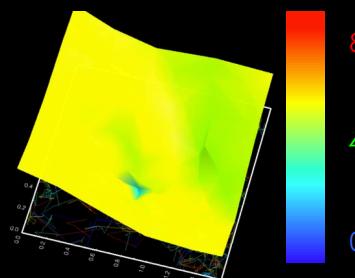
S. P. Knight *et al*, *Science* 350, 6262 p823 (2015)

Membrane Protein Interaction Energy

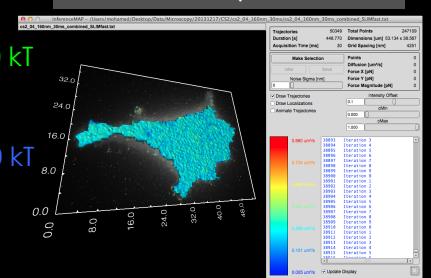


M. El Beheiry, Institut Curie

Energy of Virion Formation

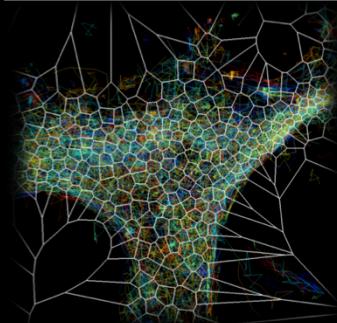


User-Friendly Interface



C. Floderer, CNRS CPBS

Adaptive Multi-Scale Meshing



# Conclusion

- Global interdisciplinary Approach
  - Physics
  - Bayesian inference
  - Statistics
- Challenges
  - Sampling complex distributions
  - Heterogeneity
  - Transition between models
  - Multi-Method Model selection
- Bayesian Inference as a global framework
  - Bacterial Chemotaxis (JBM *et al*, PNAS vol 109, No 5, p1802-1807 (2012))
  - Leucocyte search strategies 22, p2375-2382 (Sarris *et al*, Current Biology (2012))
  - Searching in turbulent environment (JBM, PNAS vol 110, No 28, p11261-11266 (2013) )

# Acknowledgments

## Pasteur Institute

- Massimo Vergassola (now at UCSD)
- Jerome Wong-Ng (now at UCSD)
- Guillaume Voisinne (now at Sloan Kettering)
- David Digregorio Lab
- Christophe Zimmer Lab
- Sven von Teeffelen Lab

## Ecole Polytechnique (LOB)

- Antigoni Alexandrou Team
- Max Richly

## Stanford

- Silvan Tuerkcan

## Colorado State University

- Diego Krapf Lab

- Institut Curie
- Maxime Dahan Lab
- Mohamed El Beheiry
- Mathieu Coppey
- Amanda Remorino
- Bassam Hajj
- Patricia Bassereau Lab

## Janelia Research Campus

- Marta Zlatic Lab
- Brian English
- Timothée Lionnet
- TIC Consortium

- Ecole Normale Supérieure
- Antoine Triller Lab
- Charlotte Salvatico
- Marianne Renner
- Christian Specht

## Grants:

ANR Grand Emprunt: Pherotaxis (2011-2015), ANR 07 NANO 062 03 (2010-2013), ANR 09 PIRI 0025 01 (2011-2014), CNANO Ile de France (2010-2012), Masson Visiting Project Janelia Research Campus (2014-2016)