# Parsimonious modeling with Information Filtering Networks:
## construction of predictive graphical models from large numbers of heterogeneous data

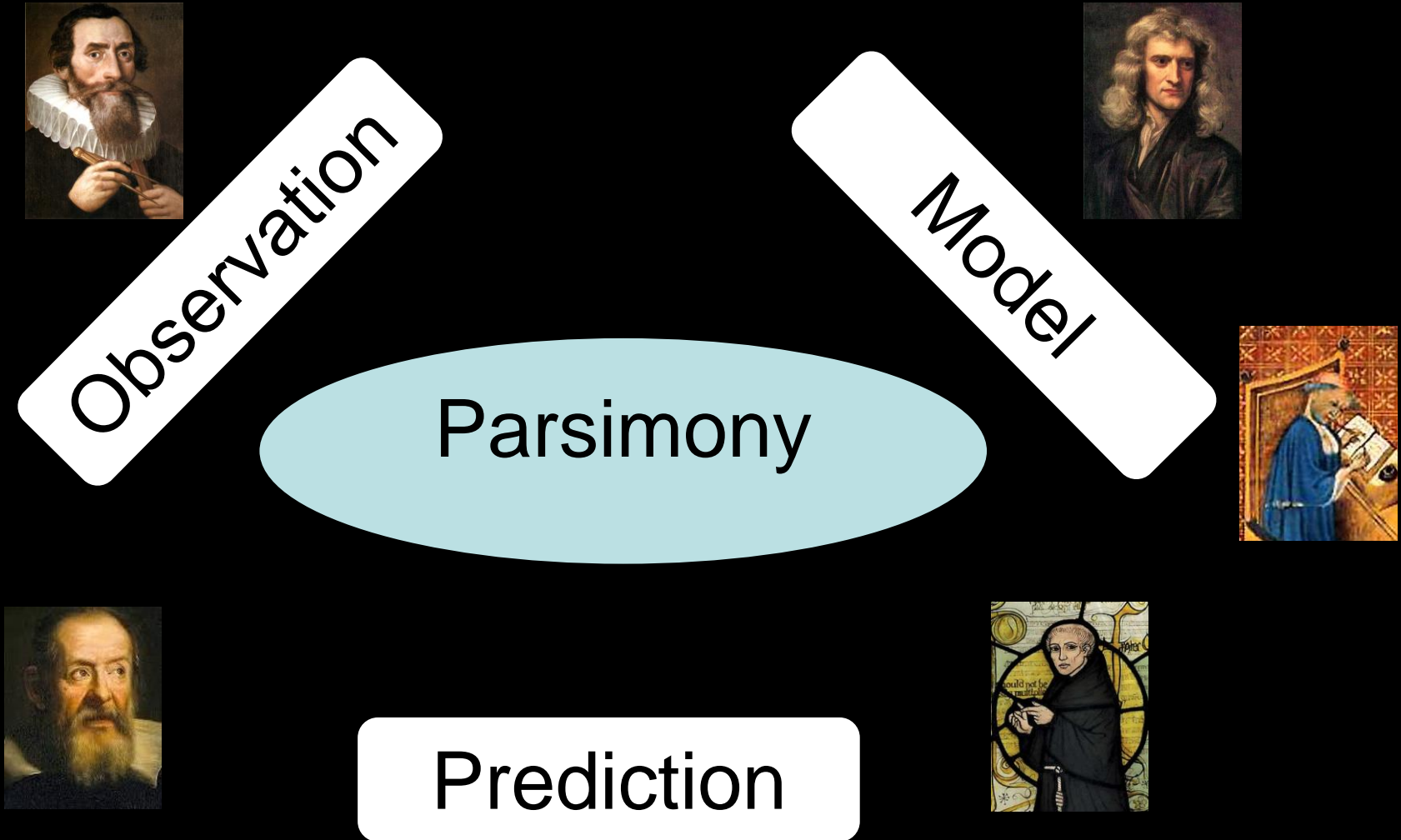*Tomaso Aste*

*&*

*Tiziana Di Matteo*

**UCL, Computer Science, Financial Computing & Analytics Group**

# Prediction

# Predictive modeling

Observation

Model

Parsimony

Prediction

Prediction is the <u>estimation of the probability</u> of a (future) event given the available information about other (past) events

$$p(X_B \mid X_A)$$

We must estimate from data the <u>most likely probability</u> distribution of the system of events

$$p(X_B \mid X_A) = \frac{p(X_A, X_B)}{p(X_A)}$$   <span style="color:red">Bayes' formula</span>

$$p(X_A, X_B)$$   <span style="color:red">joint probability</span>

**High dimensional problem!**   *(especially for big data)*

Prediction is not only about the future,

from

$$p(\mathrm{X}_B \mid \mathrm{X}_A)$$

we can <u>predict</u> the values of the variables $\mathrm{X}_B$ for any kind of scenario of the variables $\mathrm{X}_A$

**We can estimate the <u>effects</u> of events in $\mathrm{X}_A$ on $\mathrm{X}_B$**

**UCL**

The conditional probability

$$p(\mathrm{X}_B \,|\, \mathrm{X}_A)$$

is a tool for:

- **test hypothesis**
- **quantify risk**
- **stress testing**
- **analyze scenarios**

$\mathrm{X}_A$

$\mathrm{X}_B$

# Predictive modeling

Predicted future values of variables $X_B$ given past values of $X_A^-$ are the expectation values

$$E[X_B \mid X_A^-] = \sum_{X_B} X_B \, p(X_B \mid X_A^-)$$

*This is the regression and for linear models (multivariate Gaussian) this is the linear regression formula*

Uncertainty about the future given the past is quantified by the conditional entropy

$$H(X_B \mid X_A^-) = - \sum_{X_A^-, X_B} p(X_B, X_A^-) \log p(X_B \mid X_A^-)$$

The **reduction of uncertainty** on variables $X_B$ given the knowledge of the past of variables $X^-_A$ discounting for their past $X^-_B$ is

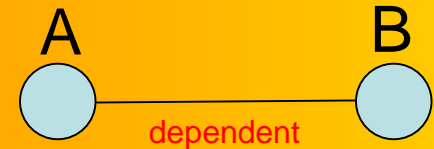$$H(X_B \mid X^-_B) - H(X_B \mid X^-_A, X^-_B) = TE(X_A \to X_B)$$

*Thisis the **tranfer entropy** that for liner models (multivariate Gaussians) coincides with **Granger causality***

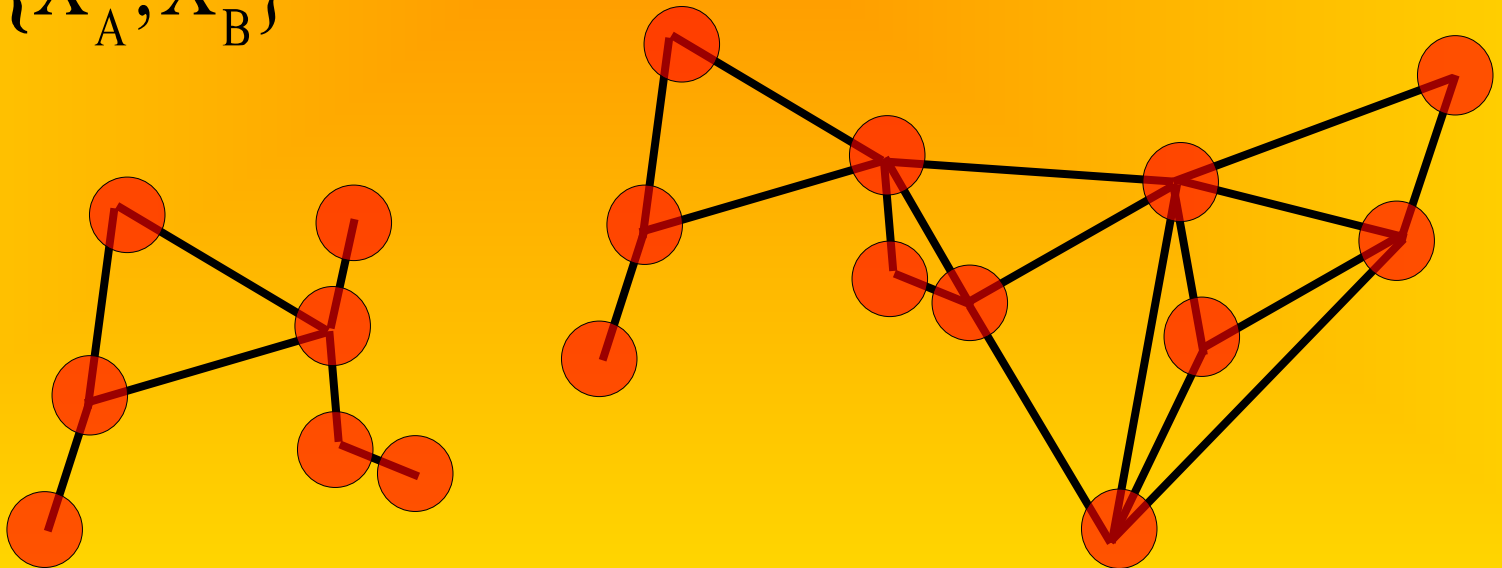To construct the joint multivariate distribution we make use of the structure of conditional dependency

$$p(X_A, X_B \mid \tilde{\mathbf{X}}) = p(X_A \mid \tilde{\mathbf{X}}) p(X_B \mid \tilde{\mathbf{X}})$$

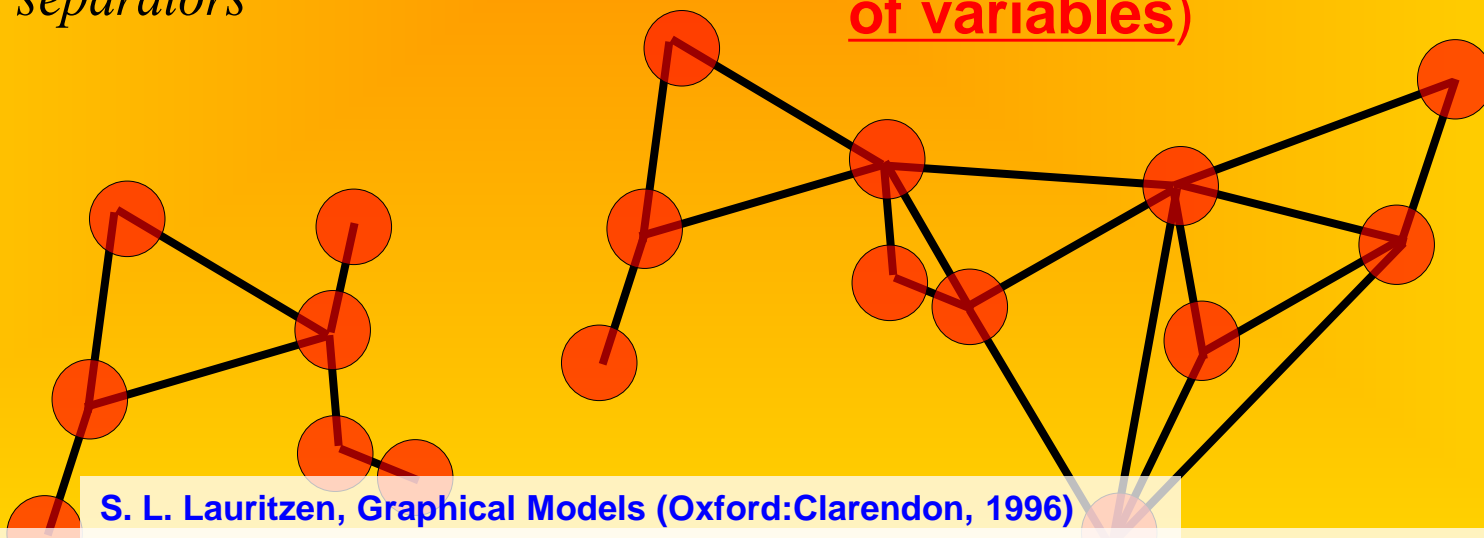$$p(X_A, X_B \mid \tilde{\mathbf{X}}) \neq p(X_A \mid \tilde{\mathbf{X}}) p(X_B \mid \tilde{\mathbf{X}})$$

$$\tilde{\mathbf{X}} = \mathbf{X} \setminus \{X_A, X_B\}$$

A      B

independent

A      B

dependent

If these inference networks are <u>chordal</u> (or decomposable) we then have

$$p(\mathbf{X}) = \frac{\tilde{\bigcirc}_{cliques} \; p(\mathbf{X}_{cliques})}{\tilde{\bigcirc}_{separators} \; p(\mathbf{X}_{separators})^{k_s - 1}}$$

The joint probability distribution of the entire system (**large number of variables**) can be estimated form the probability distributions of cliques and separators (**small number of variables**)
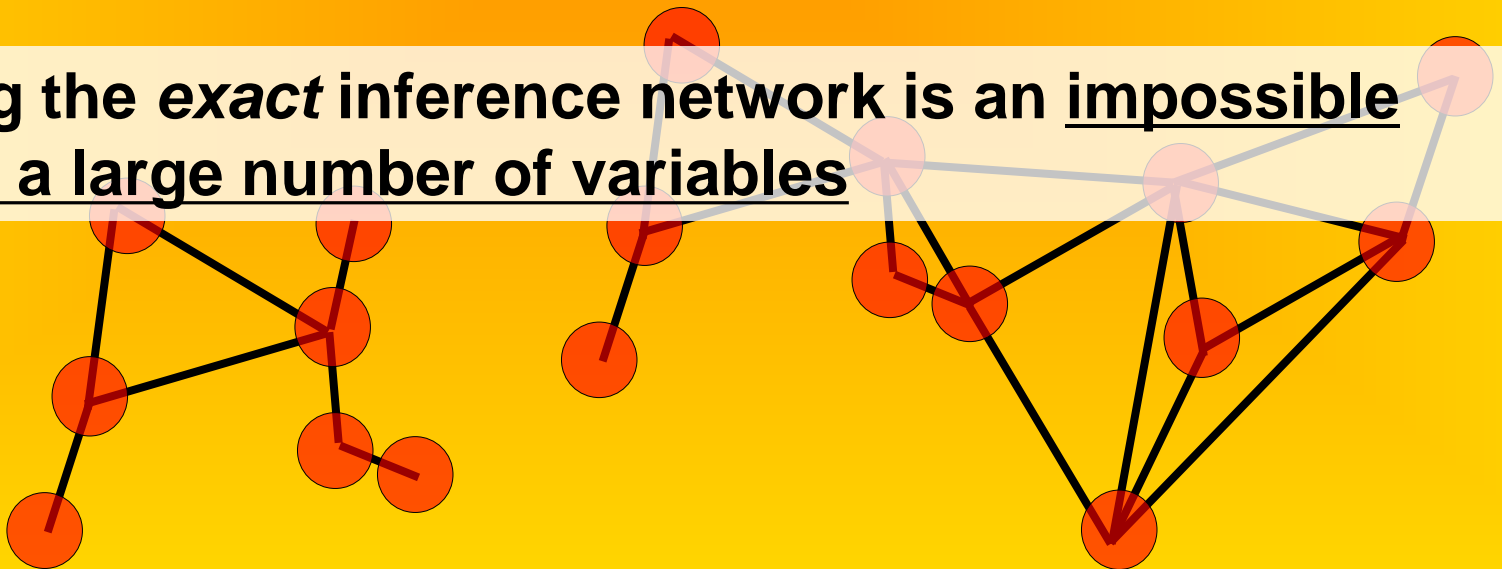
**S. L. Lauritzen, Graphical Models (Oxford:Clarendon, 1996)**
**Alexander Denev Probabilistic Graphical Models: A New Way of Thinking in Financial Modelling (Risk Books, 2015)**

This is great… however to establish conditional dependency

$$p(X_A, X_B \mid \tilde{\mathbf{X}}) \neq p(X_A \mid \tilde{\mathbf{X}}) p(X_B \mid \tilde{\mathbf{X}})$$

is very hard… actually it is <u>as hard as computing the entiere joint distribution function</u>!
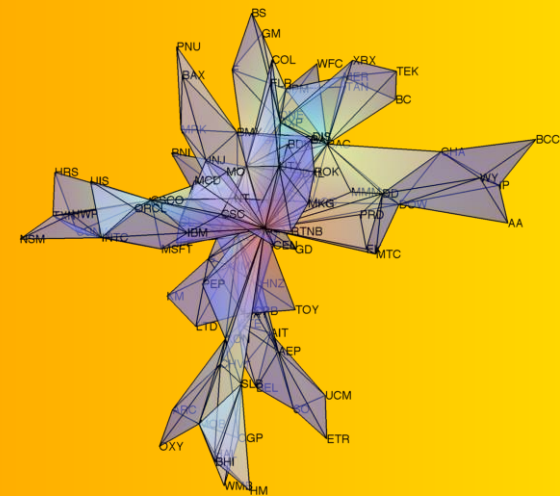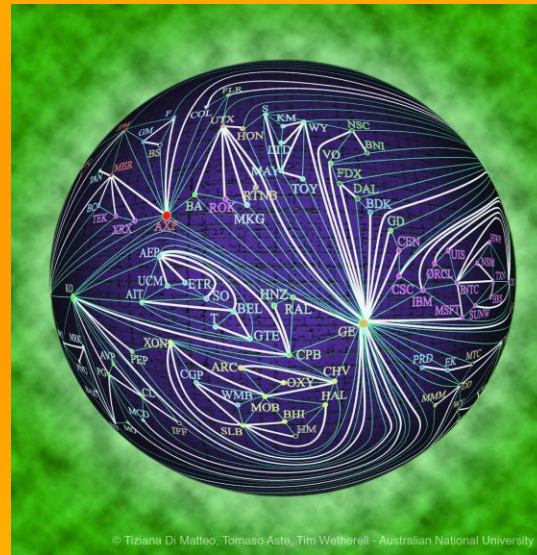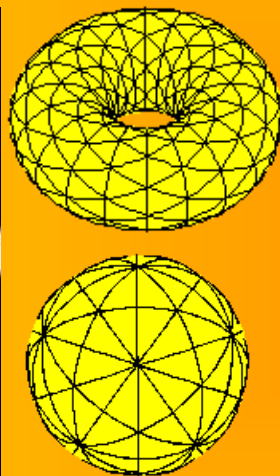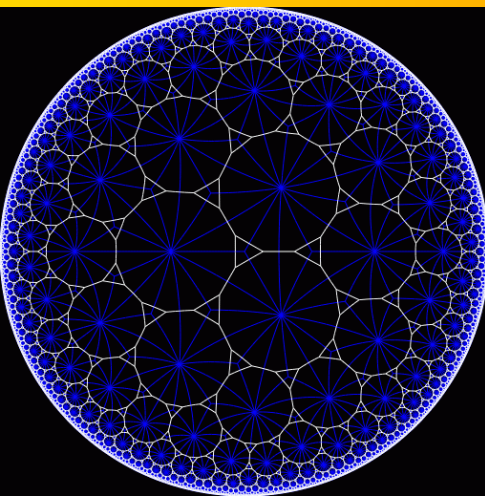
**Building the *exact* inference network is an <u>impossible</u> task for a large number of variables**

# Information filtering networks

To solve this problem we propose to build the inference structure for the graphical model as an

## Information filtering network

TA, T. Di Matteo and S. T. Hyde, *Complex networks on hyperbolic surfaces* Physica *A* 346 (2005) 20-26.



© Tiziana Di Matteo; Tomaso Aste, Tim Wetherell - Australian National University

- Massara, Guido Previde, Tiziana Di Matteo, and TA. "Network Filtering for Big Data: Triangulated Maximally Filtered Graph" Journal of Comlex Networks (2016) arXiv preprint arXiv:1505.02445 (2015).
- Nicoló Musmeci,, Tomaso Aste, and Tiziana Di Matteo. "Relation between financial market structure and the real economy: comparison between clustering methods." PloS one 10.3 (2015): e0116201.
- F. Pozzi, T. Di Matteo, and TA , "Spread of risk across financial markets: better to invest in the peripheries", Scientific Reports 3 (2013) 1665.
- W.M. Song, T. Di Matteo and T. Aste, "Hierarchical information clustering by means of topologically embedded graphs", *PLoS ONE*, 7 (2012) e31929
- M. Tumminello, T. Aste, T. Di Matteo, and R. N. Mantegna, "A tool for filtering information in complex systems" Proceedings of the National Academy of Sciences of the United States of America 102, 10421 (2005).

**Connect the nearest vertices**

*eucleadean distance = most correlated*
*hyperbolic distance = mutual information*

**Keep the graph chordal**

*clique forests*

**Add other constraints**

*max clique size (2 = MST)*
*planarity (TMFG)*
*information criteria (e.g. Akaike)*

**These are fast algorithms O(N²)**

*(topological & homological measures, betty numbers, cycles and cliques retrieved from construction)*

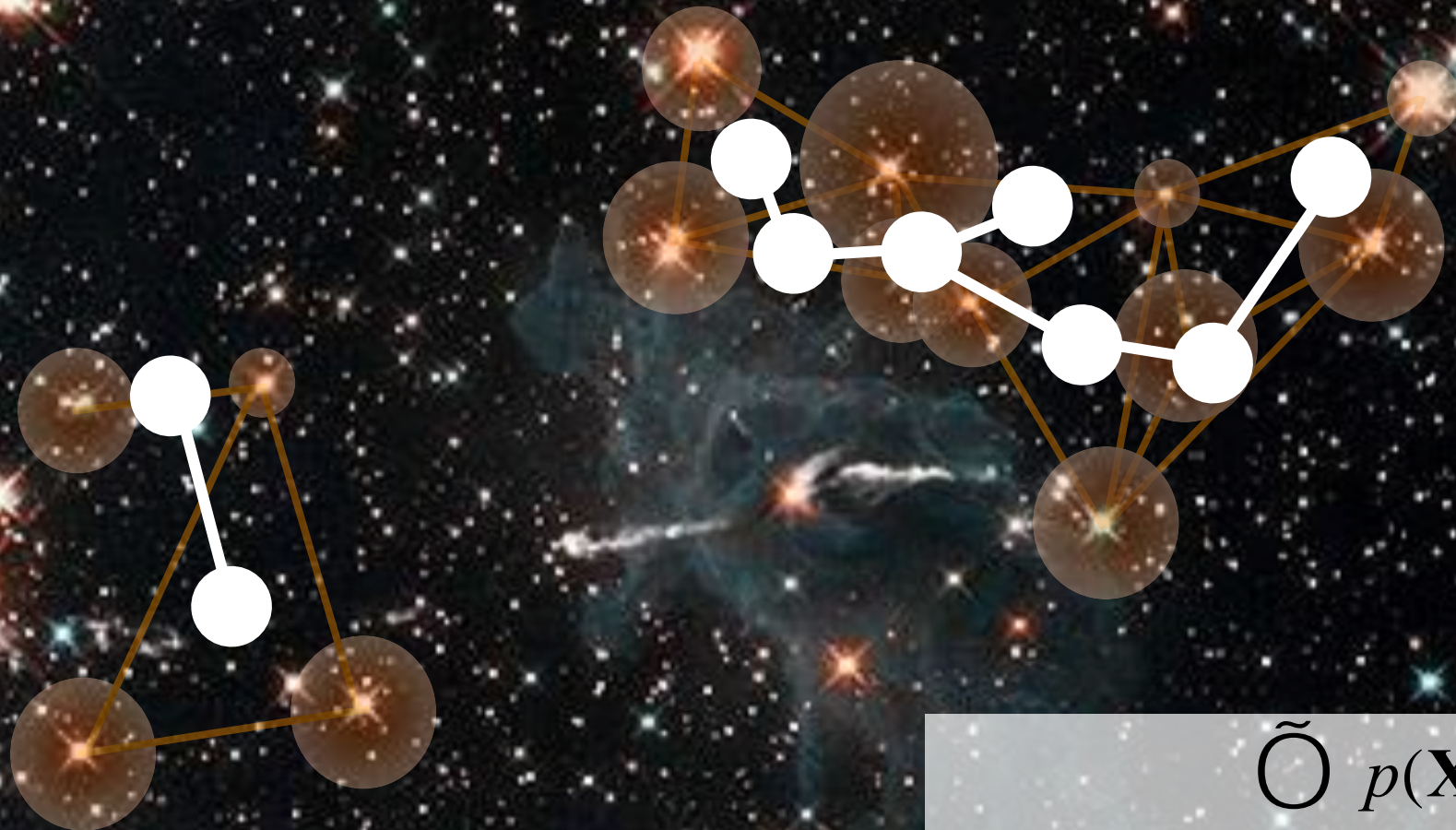Massara, Guido Previde, Tiziana Di Matteo, and TA. "Network Filtering for Big Data: Triangulated Maximally Filtered Graph" Journal of Comlex Networks (2016) arXiv preprint arXiv:1505.02445 (2015).

$$p(\mathbf{X}) = \frac{\tilde{O}_{cliques} \; p(\mathbf{X}_{cliques})}{\tilde{O}_{separators} \; p(\mathbf{X}_{separators})^{k_s - 1}}$$

$$p(\mathbf{X}) = \frac{\tilde{\prod}_{cliques} p(\mathbf{X}_{cliques})}{\tilde{\prod}_{separators} p(\mathbf{X}_{separators})^{k_s-1}}$$

$$p(\mathbf{X}) = \frac{\tilde{O}_{cliques} \; p(\mathbf{X}_{cliques})}{\tilde{O}_{separators} \; p(\mathbf{X}_{separators})^{k_s - 1}}$$

**Only low-dimensional local probabilities must be estimated**

$$p(\mathbf{X}) = \frac{\tilde{\bigcirc}_{cliques} \; p(\mathbf{X}_{cliques})}{\tilde{\bigcirc}_{separators} \; p(\mathbf{X}_{separators})^{k_s - 1}}$$

Observation

Model

Parsimony

Prediction

$$p(\mathbf{X}) = \frac{\tilde{\prod}_{cliques} p(\mathbf{X}_{cliques})}{\tilde{\prod}_{separators} p(\mathbf{X}_{separators})^{k_s-1}}$$

By constraining the model to <u>reproduce observed moments while maximizing Shannon-Gibbs entropy</u> (maximum Entropy method), at the second order, we have that the model <u>must be a multivariate Gaussian:</u>

$$p(X_1,...X_N) = \frac{1}{Z}\exp(-\sum_{i,j} X_i \mathbf{J}_{\mathbf{i,j}} X_j)$$

We keep only the significant interactions and <u>set to zero</u> (Max Ent.) the uncertain ones: $\mathbf{J}_{i,j} = 0$ iff $X_i$, $X_j$ conditionally independent

$\mathbf{J}_{i,j}$ is <u>sparse</u> and it has <u>the structure given by the information filtering network</u>

$\mathbf{J}_{i,j}$ is computed form local inversion of the covariance matrix over the clique forest

$$\mathbf{J_{i,j}} = \sum_{Cliques\ C} S(C)^{-1} - \sum_{Separators\ S} (k_S - 1)S(S)^{-1}$$

We obtain a <u>sparse inverse covariance</u> (our graphical model) by doing <u>local inversion only</u>

<u>Super-fast algorithm O(N)</u>        *even O(logN) if parallelized*

W. Barfuss, GP Massara, T Di Matteo & TA "Parsimonious modeling with Information Filtering Networks" arXiv preprint arXiv:1602.07349 (2016).

# Test:

is our model, built form past observations, associated with a large likelihood for future observations?

**Observation**

P a s t

**Model**

**Prediction & Testing**

$$p(System) = \frac{\widetilde{\prod_{cliques}} p(cliques)}{\widetilde{\prod_{separators}} p(separators)^{k_s - 1}}$$
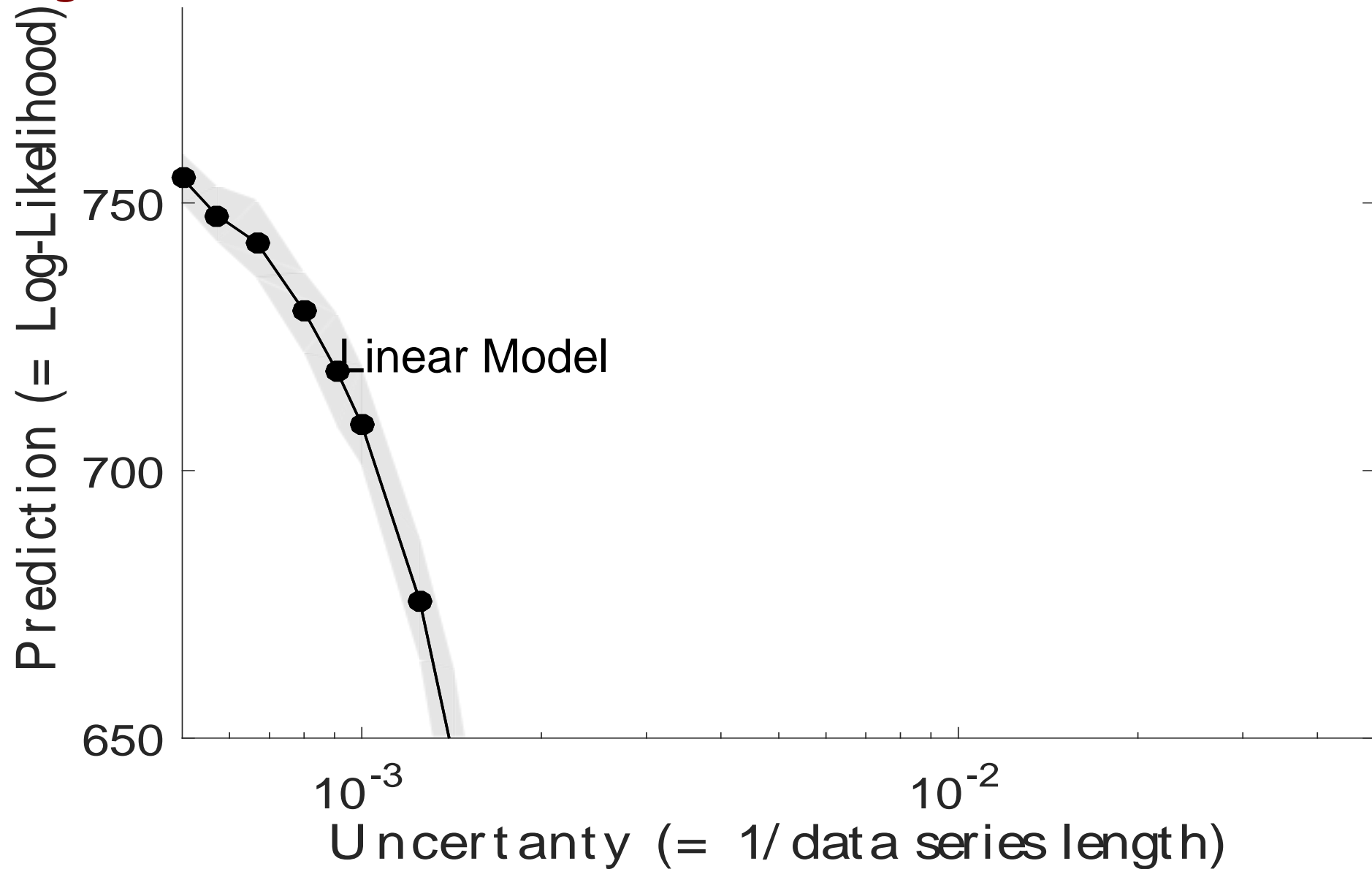
Statistical description

F u t u r e

## is our model, built form past observations, associated with a large likelihood for future observations?



Linear Model

# Test:

**is our model, built form past observations, associated with a large likelihood for future observations?**



Prediction (= Log-Likelihood)

750

700

650

Linear Model

State-of-the-art sparse model (G-lasso)

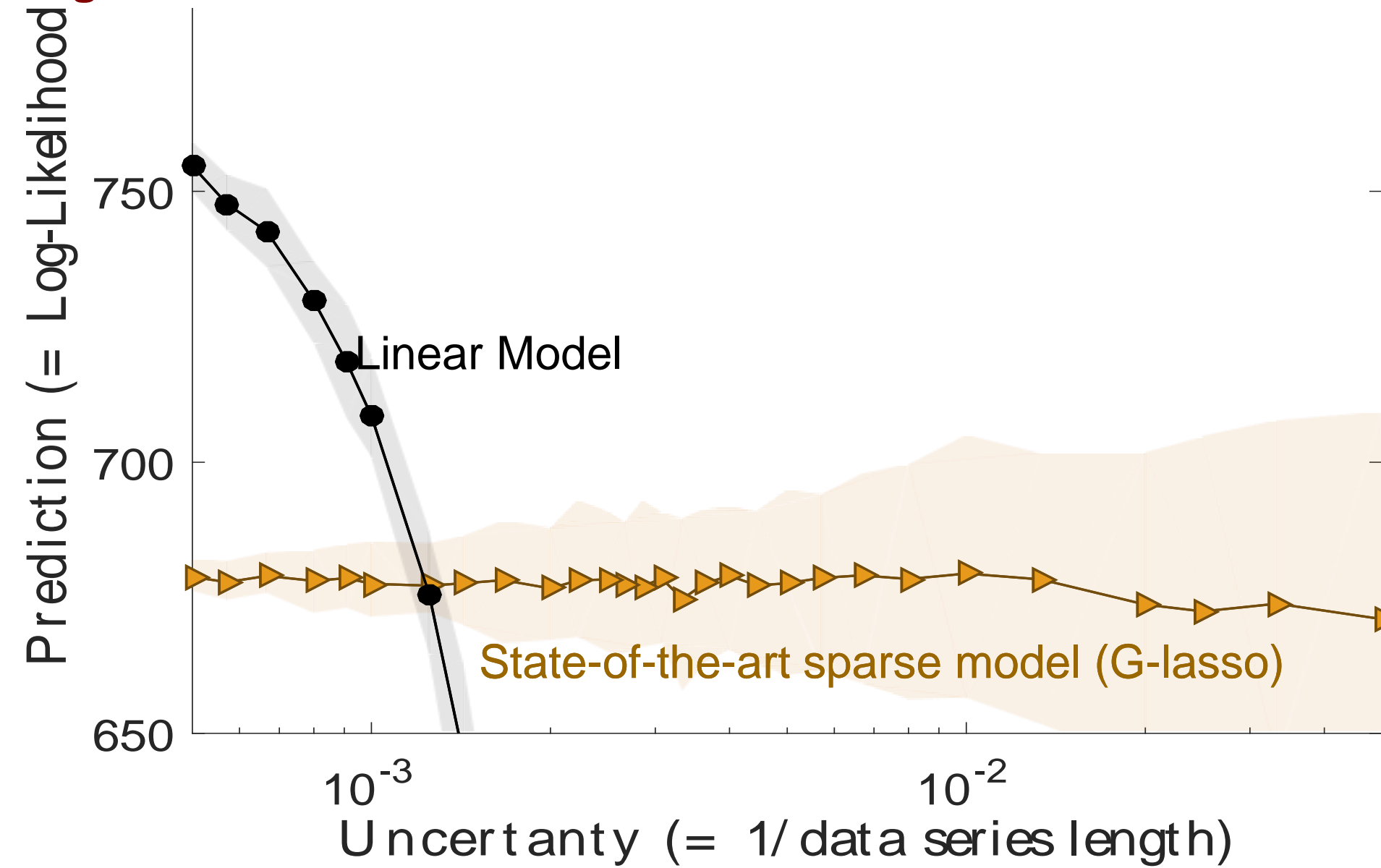$10^{-3}$      $10^{-2}$

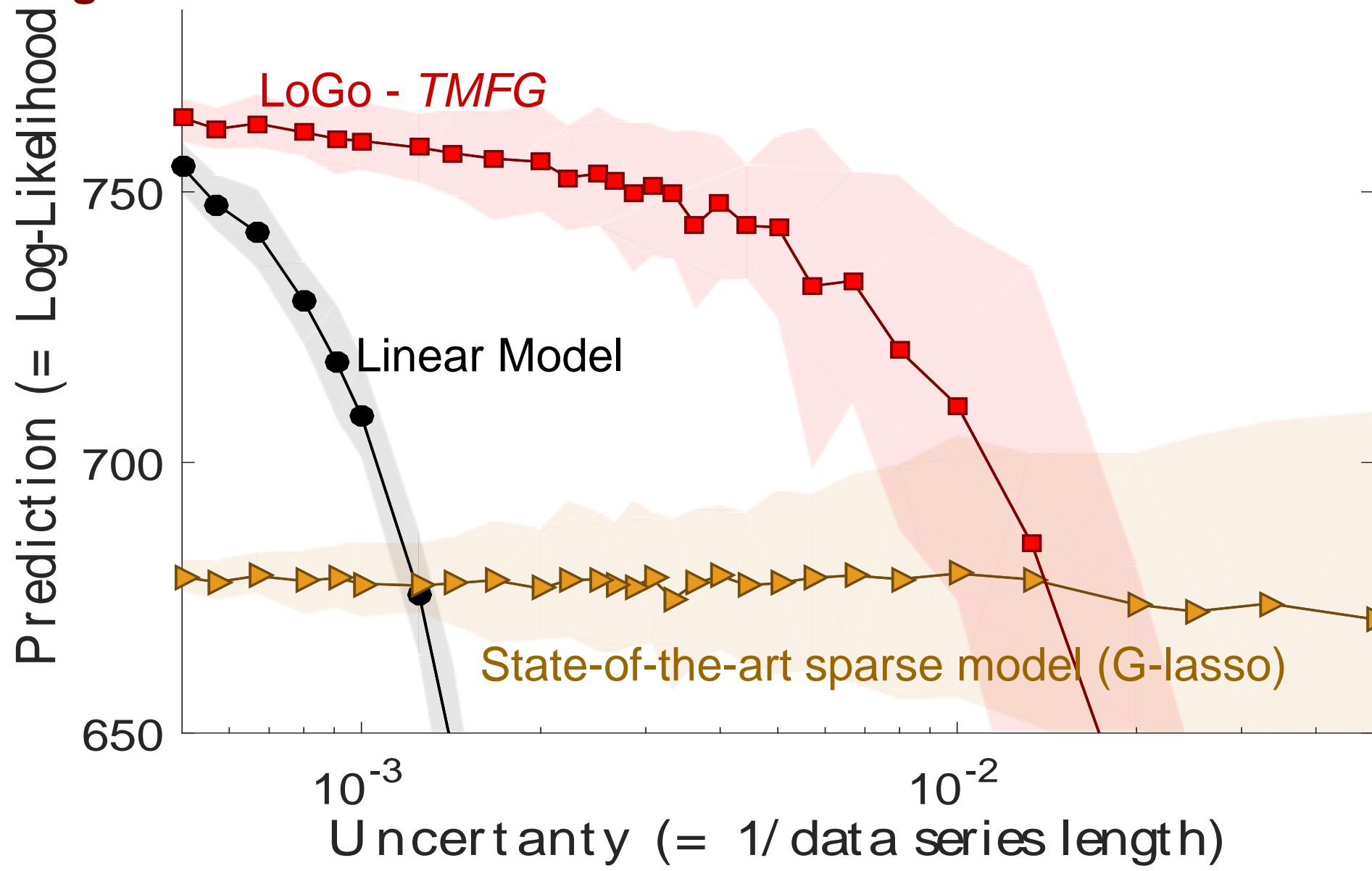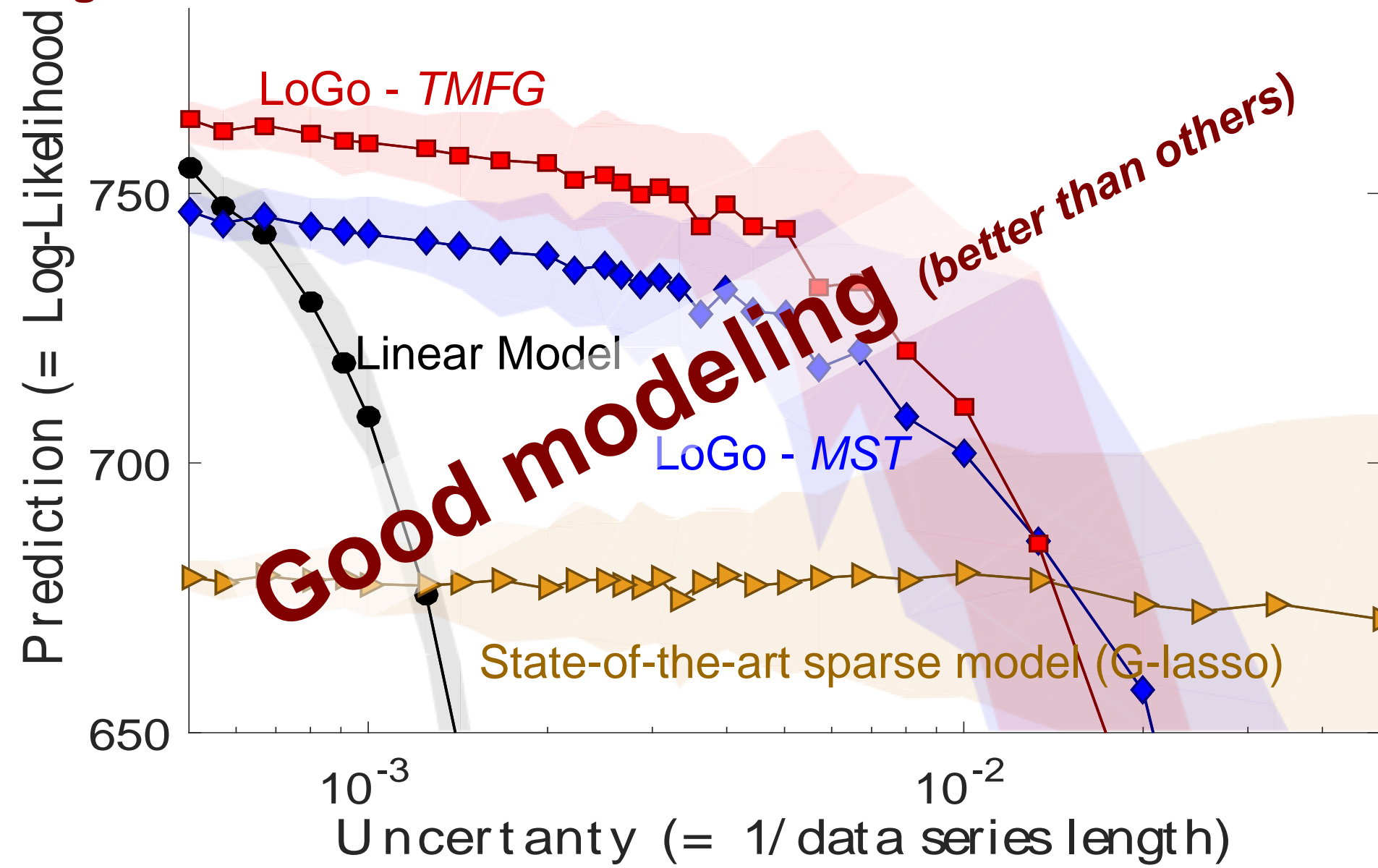Uncertanty (= 1/ data series length)

Test:
is our model, built form past observations, associated with a large likelihood for future observations?

Test:

is our model, built form past observations, associated with a large likelihood for future observations?

In which sense we predict?

**With p(X$_{future}$|X$_{past}$) we can predict the future**

This is the same as (linear) <u>regression</u>

$$E[X_B \mid X_A] = \mathring{\sum_{X_B}} X_B p(X_B \mid X_A) = -\mathbf{J}_{\mathbf{BB}}^{-1}\mathbf{J}_{\mathbf{BA}}X_A$$

and also <u>Granger causality (2x Transfer entropy)</u>

$$G(X_A \rightarrow X_B) = \log|\mathbf{J}_{\mathbf{B^-B^-}}| - \log|\mathbf{J}_{\mathbf{B^-B^-}} - \mathbf{J}_{\mathbf{B^-A^-}}\mathbf{J}_{\mathbf{A^-A^-}}^{-1}\mathbf{J}_{\mathbf{A^-B^-}}|$$

The advantage is that we have a **<u>sparse model</u>** computed in a very efficient way applicable to big-data predictive analytics

# Test:

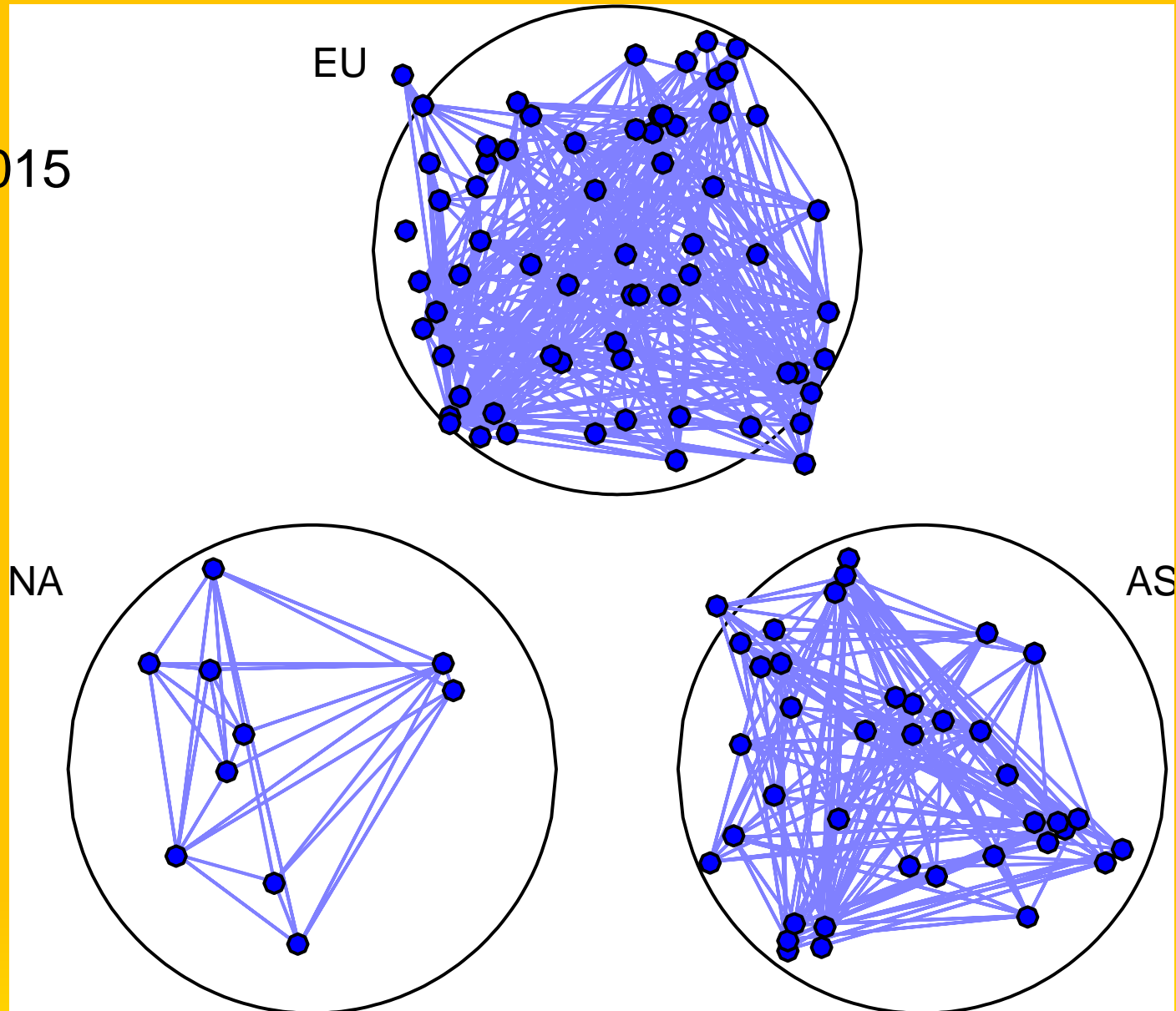## Uncertainty spillover across regions in banking system
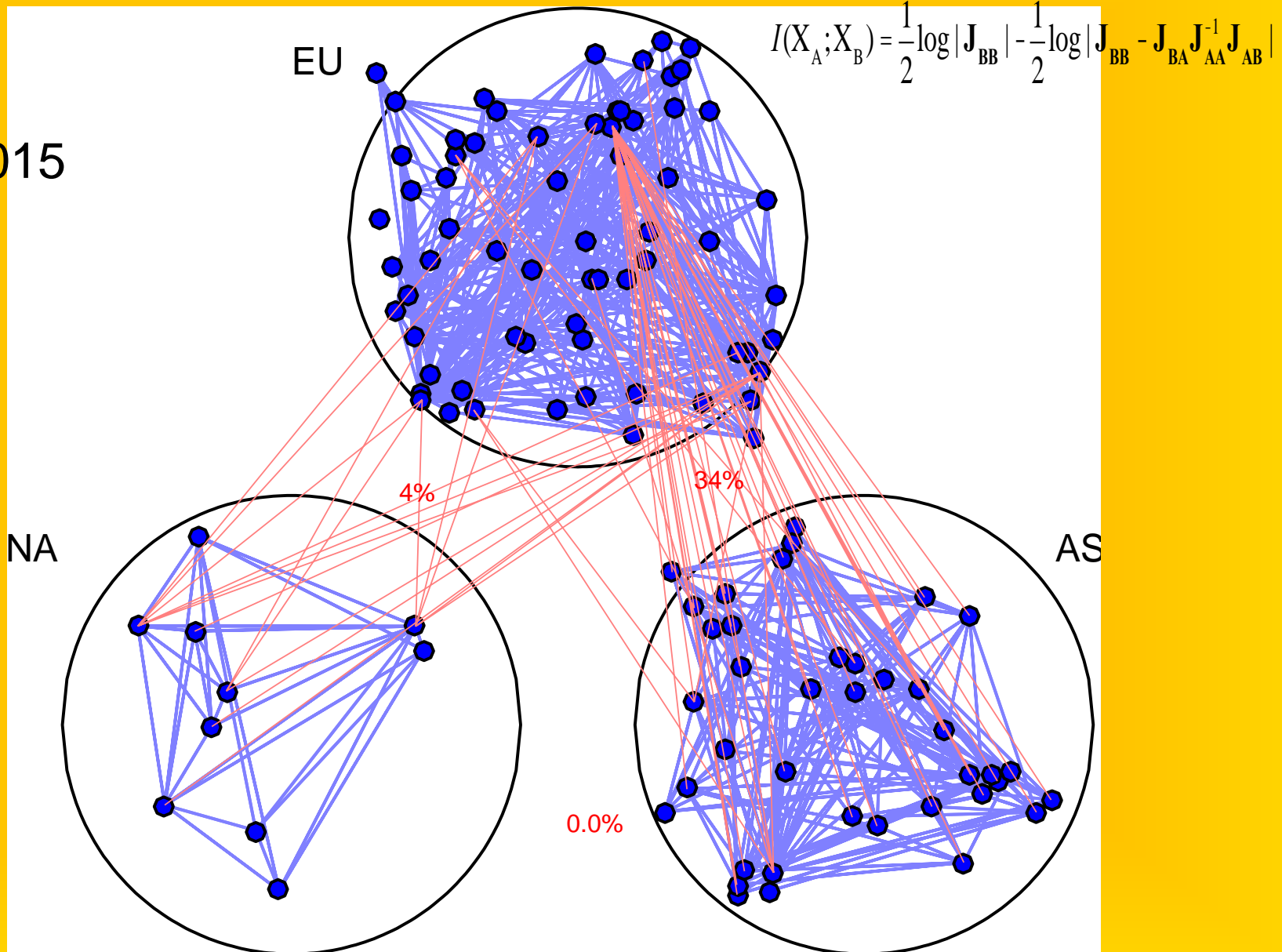
**2005-2015**

115 banks:
10 in NA , 66 in EU, 39 in AS

## Uncertainty spillover across regions in banking system

2005-2015

# Test:

## Uncertainty spillover across regions in banking system

2005-2015



$$I(X_A; X_B) = \frac{1}{2} \log |\mathbf{J}_{BB}| - \frac{1}{2} \log |\mathbf{J}_{BB} - \mathbf{J}_{BA} \mathbf{J}_{AA}^{-1} \mathbf{J}_{AB}|$$

EU

NA

AS

4%

34%

0.0%

# Test:

## Uncertainty spillover across regions in banking system



2005-2008

$$I(X_A; X_B) = \frac{1}{2}\log|\mathbf{J}_{BB}| - \frac{1}{2}\log|\mathbf{J}_{BB} - \mathbf{J}_{BA}\mathbf{J}_{AA}^{-1}\mathbf{J}_{AB}|$$

$$TE(X_A \rightarrow X_B) = \frac{1}{2}\log|\mathbf{J}_{B^-B^-}| - \frac{1}{2}\log|\mathbf{J}_{B^-B^-} - \mathbf{J}_{B^-A^-}\mathbf{J}_{A^-A^-}^{-1}\mathbf{J}_{A^-B^-}|$$

EU

NA

AS

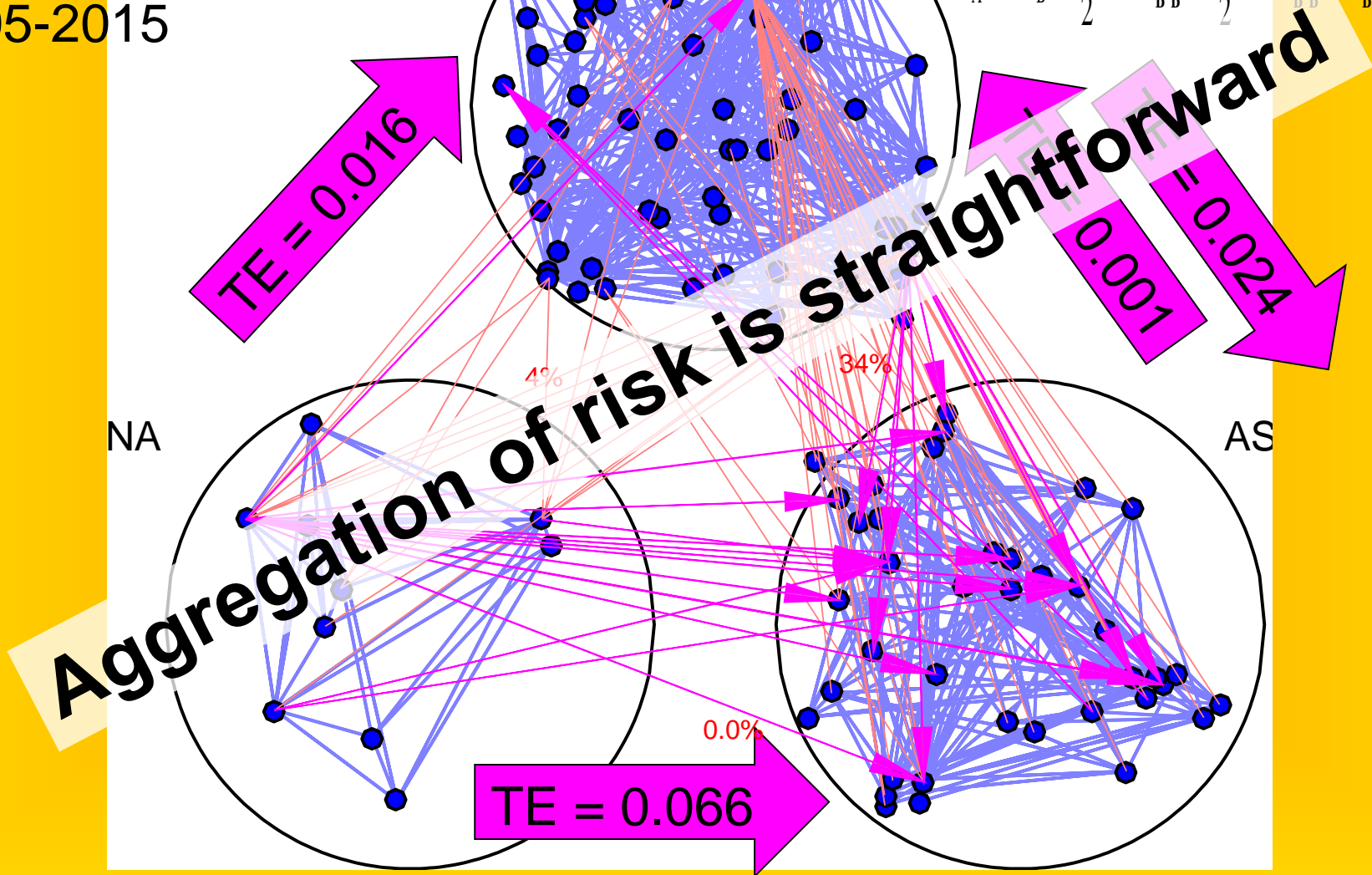TE = 0.014

TE = 0.150

TE = 0.057

TE = 0.048

8%

23%

# Test:
## Uncertainty spillover across regions in banking system

2008-2012
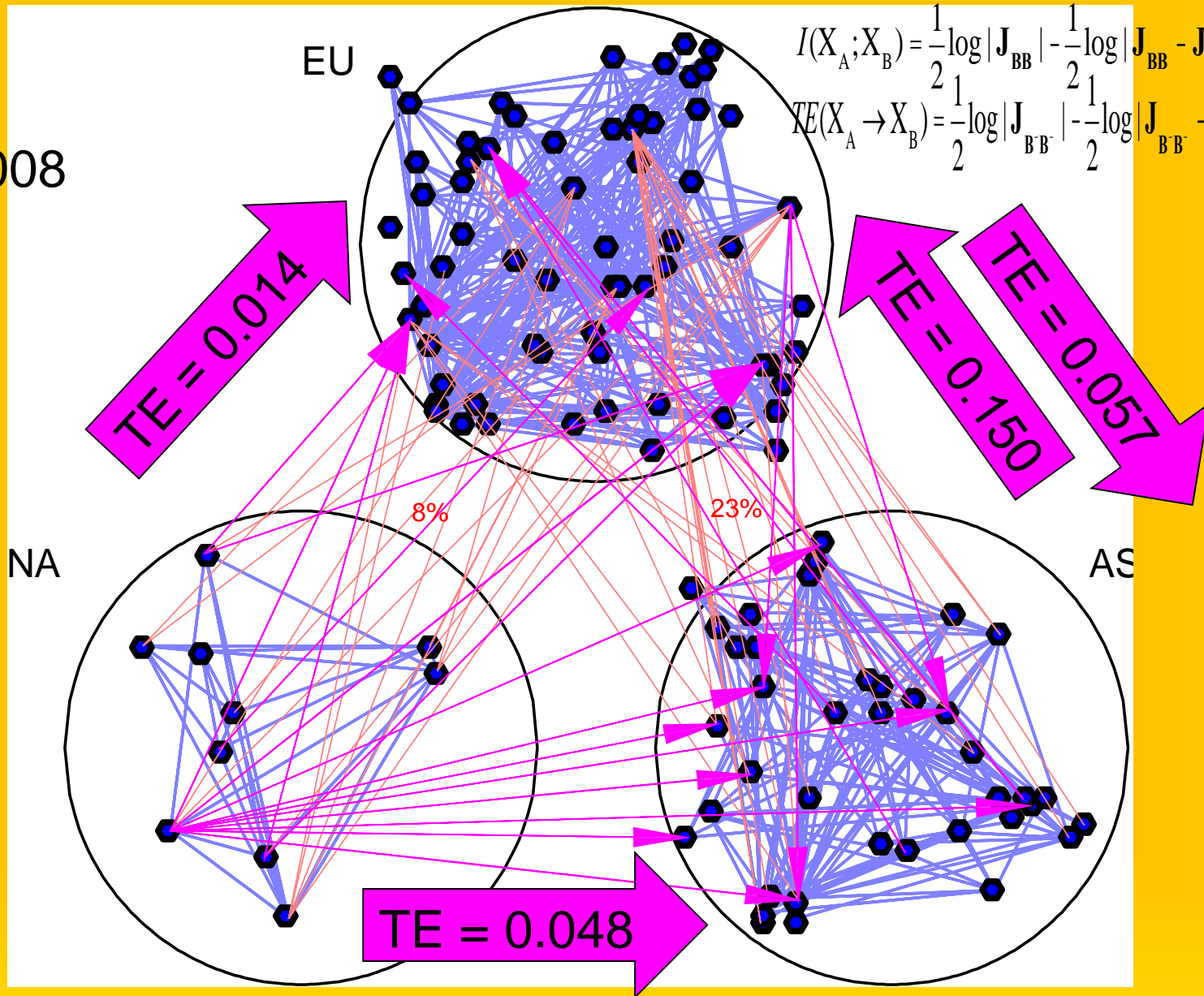
$$I(X_A; X_B) = \frac{1}{2}\log|\mathbf{J}_{BB}| - \frac{1}{2}\log|\mathbf{J}_{BB} - \mathbf{J}_{BA}\mathbf{J}_{AA}^{-1}\mathbf{J}_{AB}|$$

$$TE(X_A \to X_B) = \frac{1}{2}\log|\mathbf{J}_{B^-B^-}| - \frac{1}{2}\log|\mathbf{J}_{B^-B^-} - \mathbf{J}_{B^-A^-}\mathbf{J}_{A^-A^-}^{-1}\mathbf{J}_{A^-B^-}|$$

EU

NA

AS

TE = 0.0

TE = 0.0

TE = 0.0

TE = 0.0

16%

42%

# Test:

## Uncertainty spillover across regions in banking system

2012-2015

$$I(X_A; X_B) = \frac{1}{2}\log|\mathbf{J}_{BB}| - \frac{1}{2}\log|\mathbf{J}_{BB} - \mathbf{J}_{BA}\mathbf{J}_{AA}^{-1}\mathbf{J}_{AB}|$$

$$TE(X_A \rightarrow X_B) = \frac{1}{2}\log|\mathbf{J}_{B^-B^-}| - \frac{1}{2}\log|\mathbf{J}_{B^-B^-} - \mathbf{J}_{B^-A^-}\mathbf{J}_{A^-A^-}^{-1}\mathbf{J}_{A^-B^-}|$$

EU

NA

AS

TE = 0.001

TE = 0.01

TE = 0.009

TE = 0.002

2%

25%

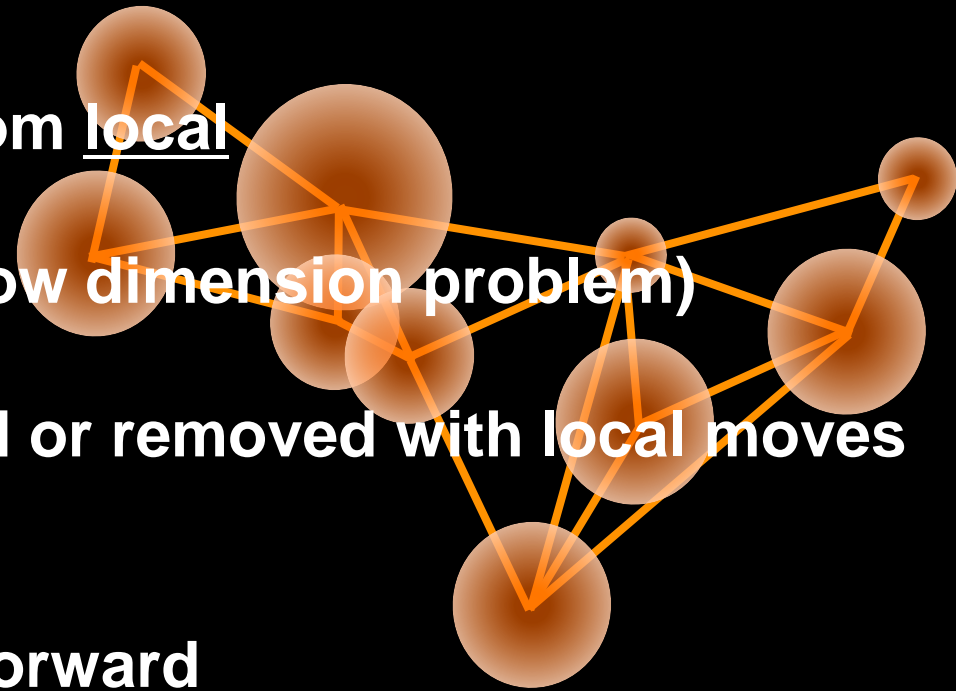With $p(X_B|X_A-)$ we can quantify probability of future events

With $p(X_B|X_A)$ we can predict impact of unobserved scenarios and test hypothesis

$p(X_A,X_B)$ can be constructed from <u>local</u> probability estimations over an <u>information filtering network</u> (low dimension problem)

Nodes and edges can be added or removed with local moves only

Aggregation of risk is straightforward

LoGo works better than state-of-the-art sparse graphical models and it is faster

# Thank YOU!

W. Barfuss, GP Massara, T Di Matteo & TA "Parsimonious modeling with Information Filtering Networks" arXiv preprint arXiv:1602.07349 (2016).

Massara, Guido Previde, Tiziana Di Matteo, and TA. "Network Filtering for Big Data: Triangulated Maximally Filtered Graph" Journal of Comlex Networks (2016) arXiv preprint arXiv:1505.02445 (2015).

http://www.cs.ucl.ac.uk/staff/tomaso_aste/

http://fincomp.cs.ucl.ac.uk/

http://blockchain.cs.ucl.ac.uk/

*Si l'ordre satisfait la raison, le désordre fait les délices de l'imagination*
*Paul Claudel*