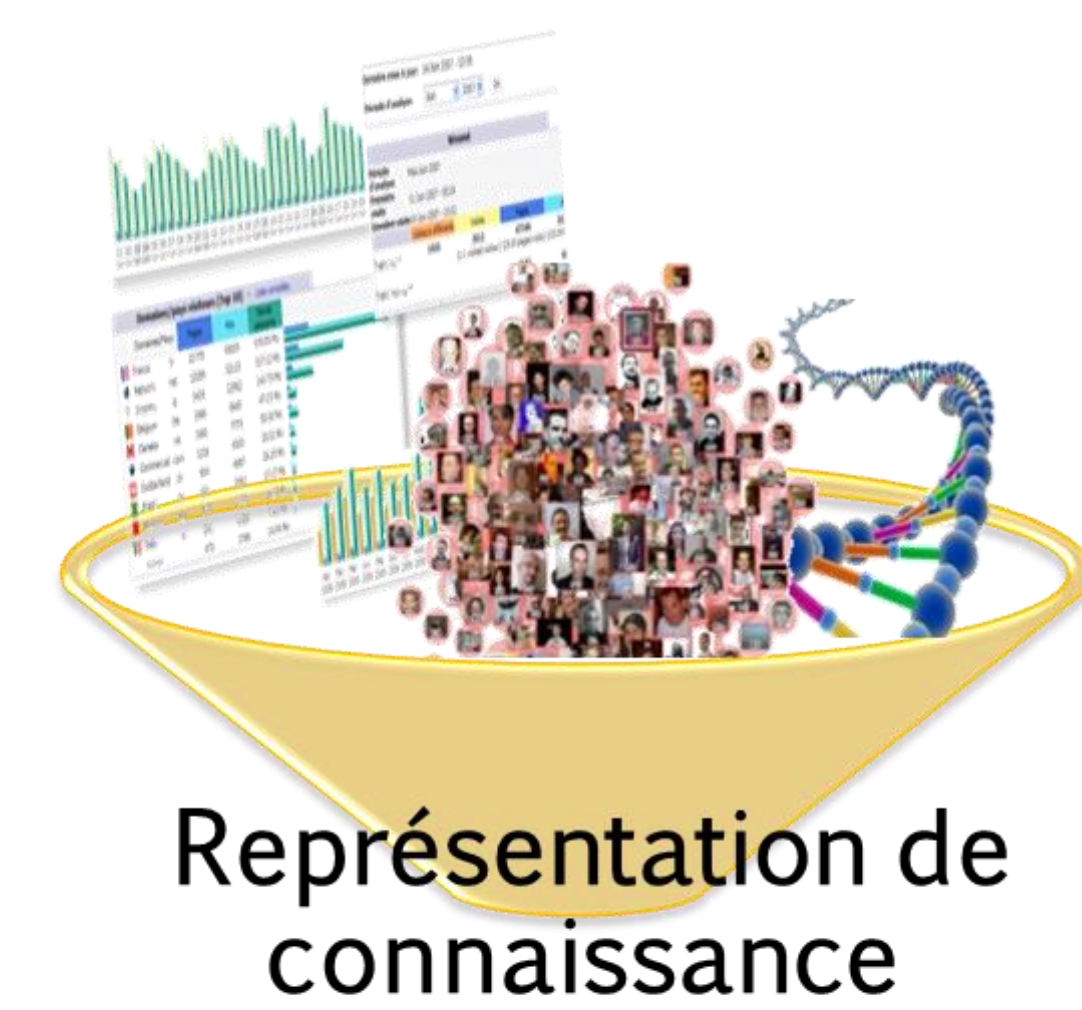


3. GRANDES MASSES DE DONNEES



Représentation de connaissance

Parties prenantes



LGI2P – équipe KID

Equipe

- Sylvie Ranwez
- Jacky Montmain
- Michel Crampes
- Gérard Dray
- Stefan Janaqi
- Michel Plantié
- François Troussel

- Nicolas Fiorini
- Sébastien Harispe

Partenaires



UNIVERSITAT ROVIRA I VIRGILI



Vincent Ranwez¹

Montserrat Batet²

David Sánchez²

¹ Equipe DAVEM, SupAgro Montpellier,

² Université Rovira i Virgili de Tarragone, Espagne



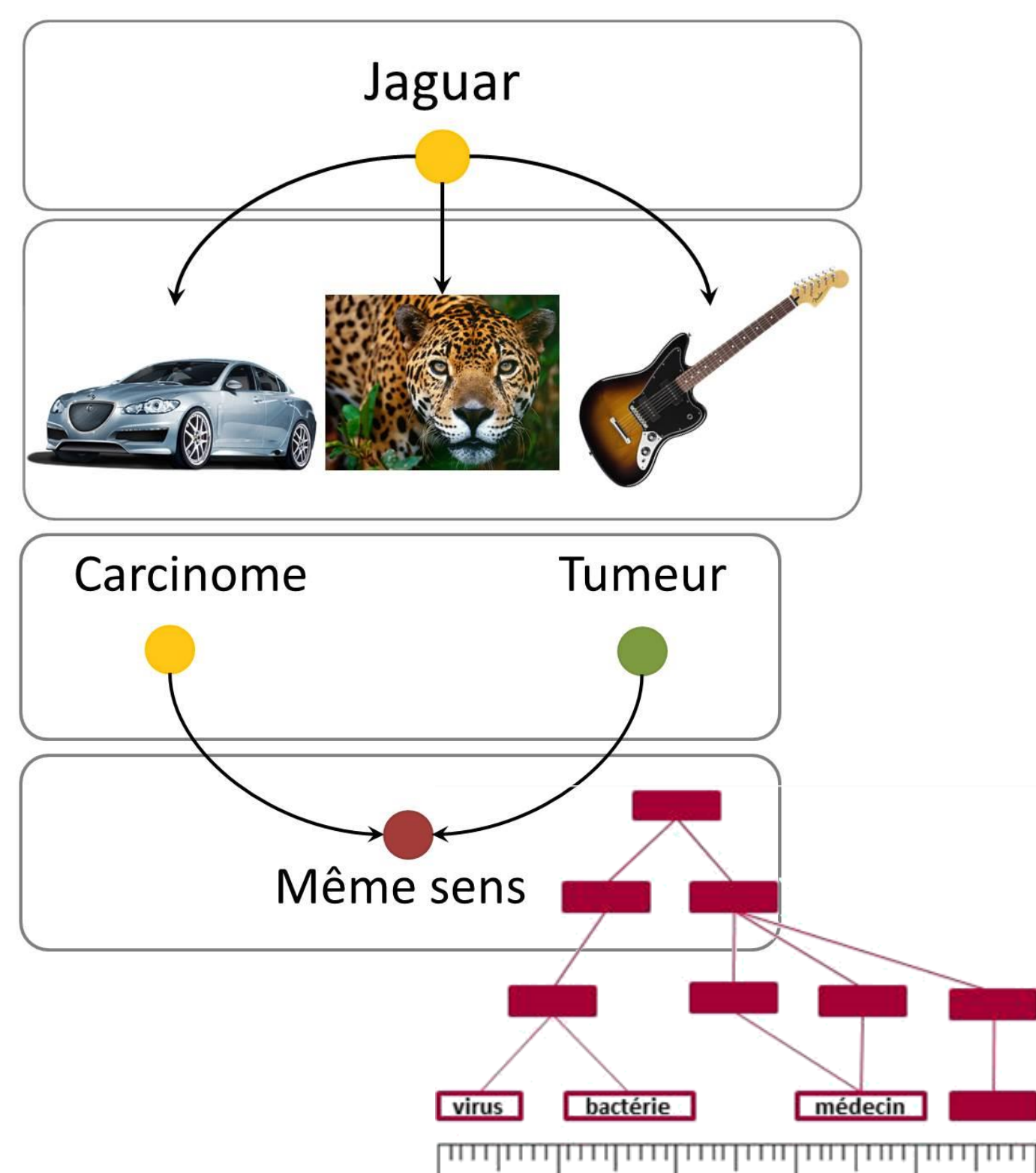
Sources d'information diverses, hétérogènes

Données non structurées

- Corpus de **textes** (publications, pages Web, réseaux sociaux)
- Ressources **multimédia** (images, sons, vidéos...)

Données structurées

- Données **liées**, **massives** et à **caractère sémantique**
- **Représentation de la connaissance** (ontologies), graphes RDF(S), OWL,...



Recherche d'information conceptuelle

Améliorer la pertinence des résultats

- Désambigüiser, généraliser, spécialiser les requêtes en utilisant une **ontologie** de domaine
- Assurer de meilleurs taux de précision et de rappel par des **mesures sémantiques** appropriées

Personnalisation, visualisation et interactivité

- Paramétrisation et **personnalisation** du système et de l'interface : pondération, reformulation, lentilles
- Visualisation globale des résultats sur une **carte sémantique** : affichage en 2D des résultats en fonction de leur degré de pertinence
- Justification des résultats

SML – Semantic Measures Library

Librairie logicielle dédiée au calcul de **similarité sémantique**

- Open-source, langage Java, exécutable en ligne de commande
- Traitement de gros volumes de données

OBIRS – Ontology Based Information Retrieval System

Recherche d'information conceptuelle centrée sur l'utilisateur

- Application à la **recherche de gènes** indexés par la Gene Ontology
- Recherche de **publications scientifiques biomédicales**

Système de recommandation

Exploitation des mesures de similarité dans le contexte des données liées pour un système de recommandation

Kalitmo

Visualisation de données structurées pour la gestion de collectifs

Développements logiciels

Extraction d'opinion

- Fouille de **textes** et **apprentissage** supervisé sur de larges corpus (Web) : TAL, **segmentation**, analyses statistiques
- Identification d'un **lexique**, détection de critères relatifs à un domaine
- **Polarisation contextuelle**, extraction conceptuelle et évaluation de critères

Synopsis

Outil de détection d'opinion, prise en compte d'une évolution temporelle

CoLexIR

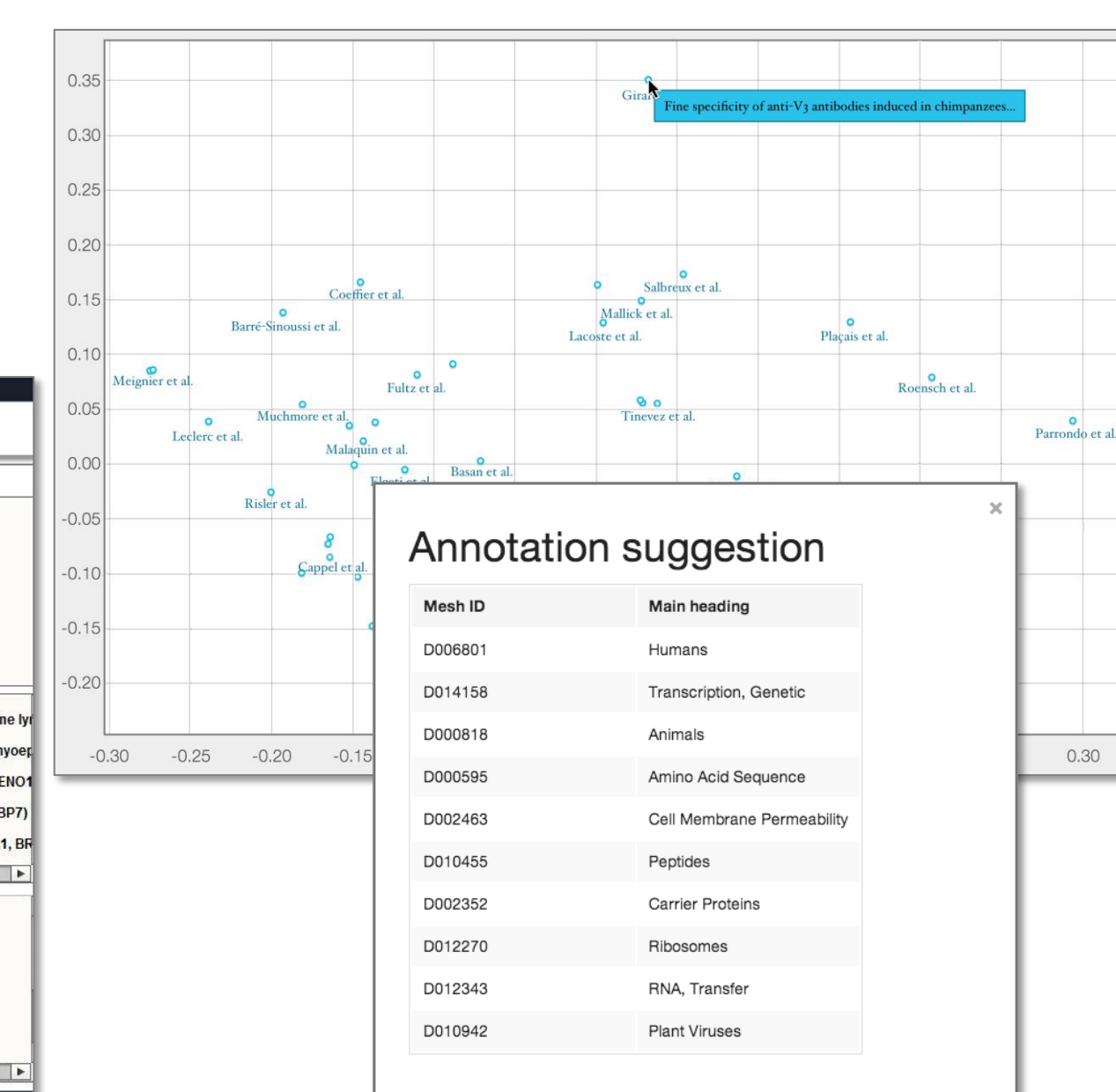
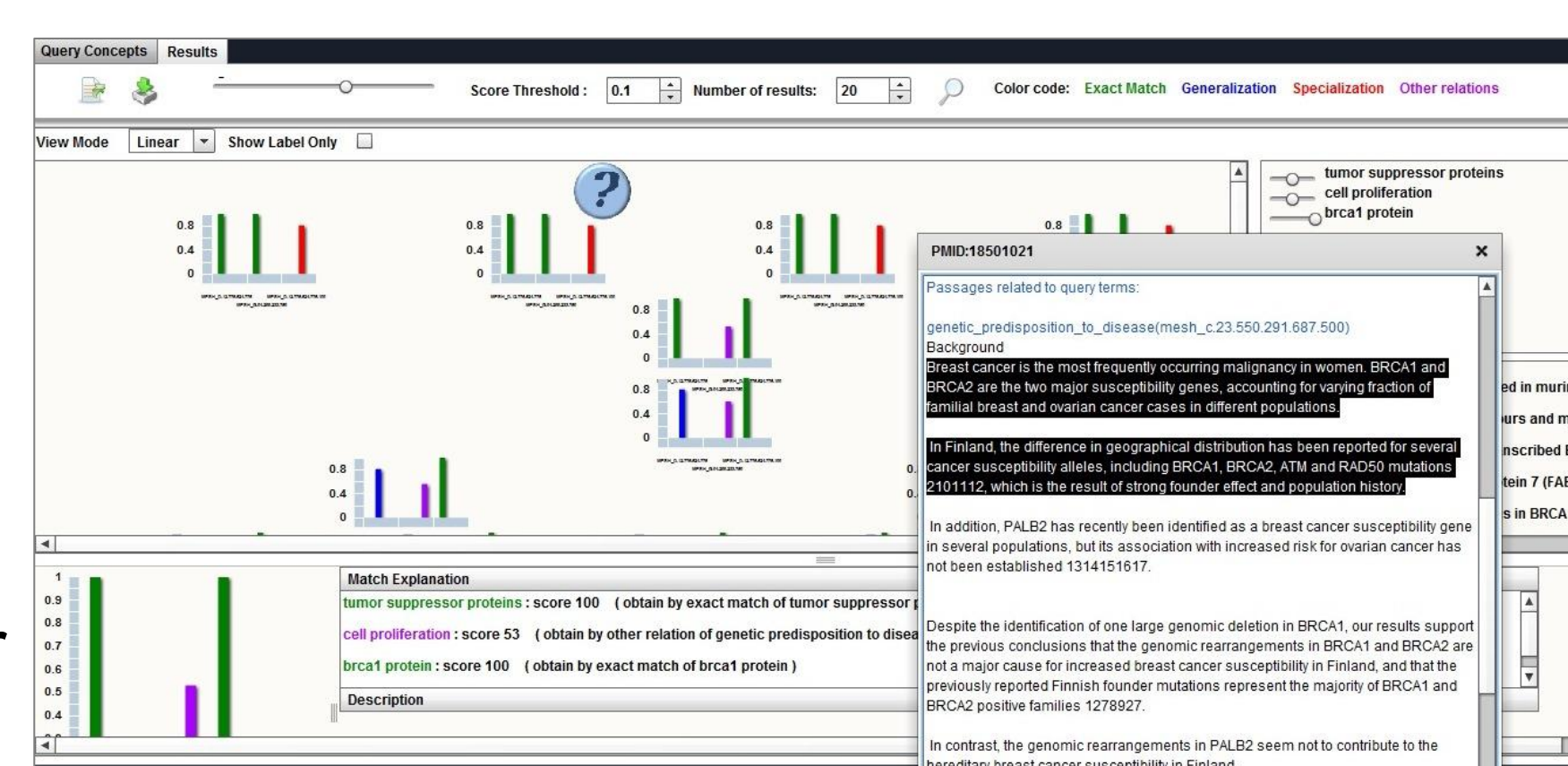
Environnement hybride de RI associant deux approches conceptuelle et lexicale. Application à la fouille de corpus scientifiques avec mise en évidence des passages pertinents pour l'utilisateur

Développements logiciels

Indexation conceptuelle

Indexation par propagation

- Interface interactive d'assistance à l'indexation
- Annotation conceptuelle semi-automatique



Parties prenantes

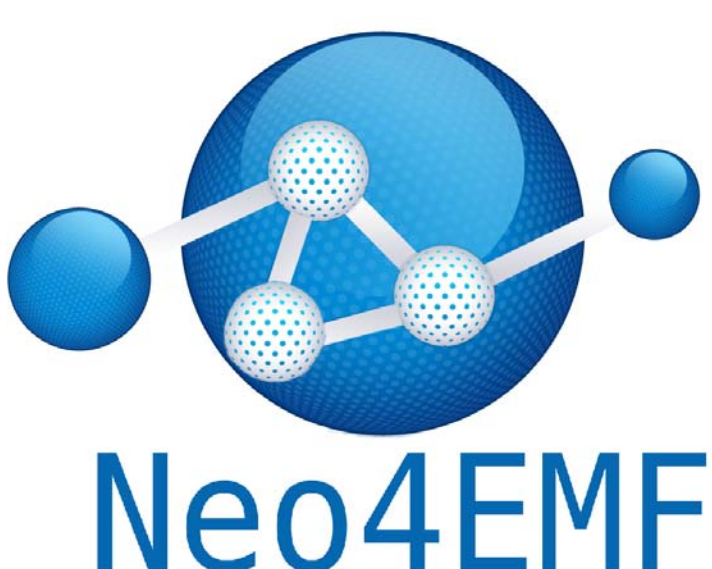



Auteurs

Benelallam Amine¹
Tisi Massimo¹
Sunyé Gerson²
Gomez Llana Abel¹

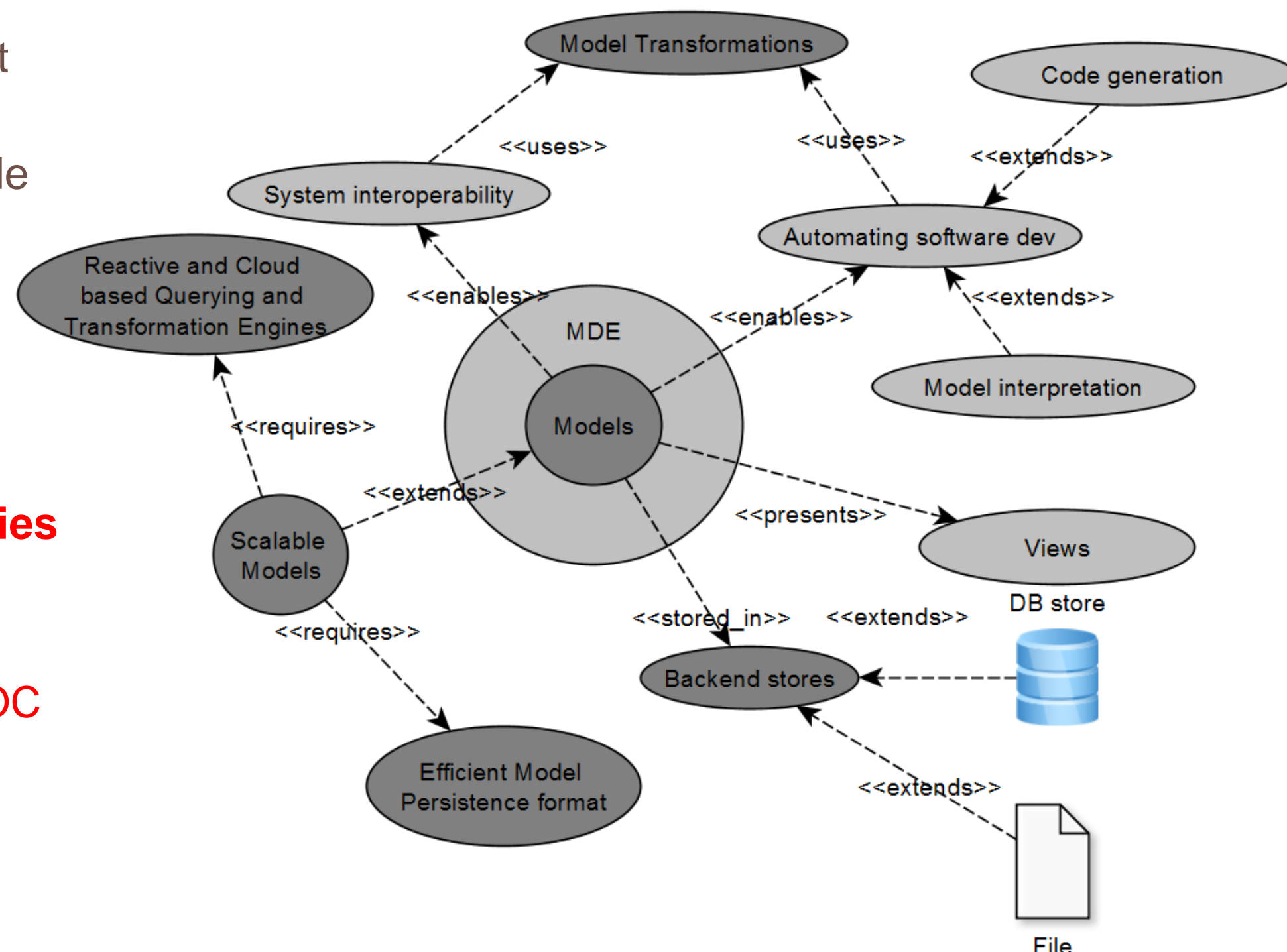
¹ AtlanMod, Mines Nantes, INRIA, LINA

² AtlanMod, Université de Nantes, LINA

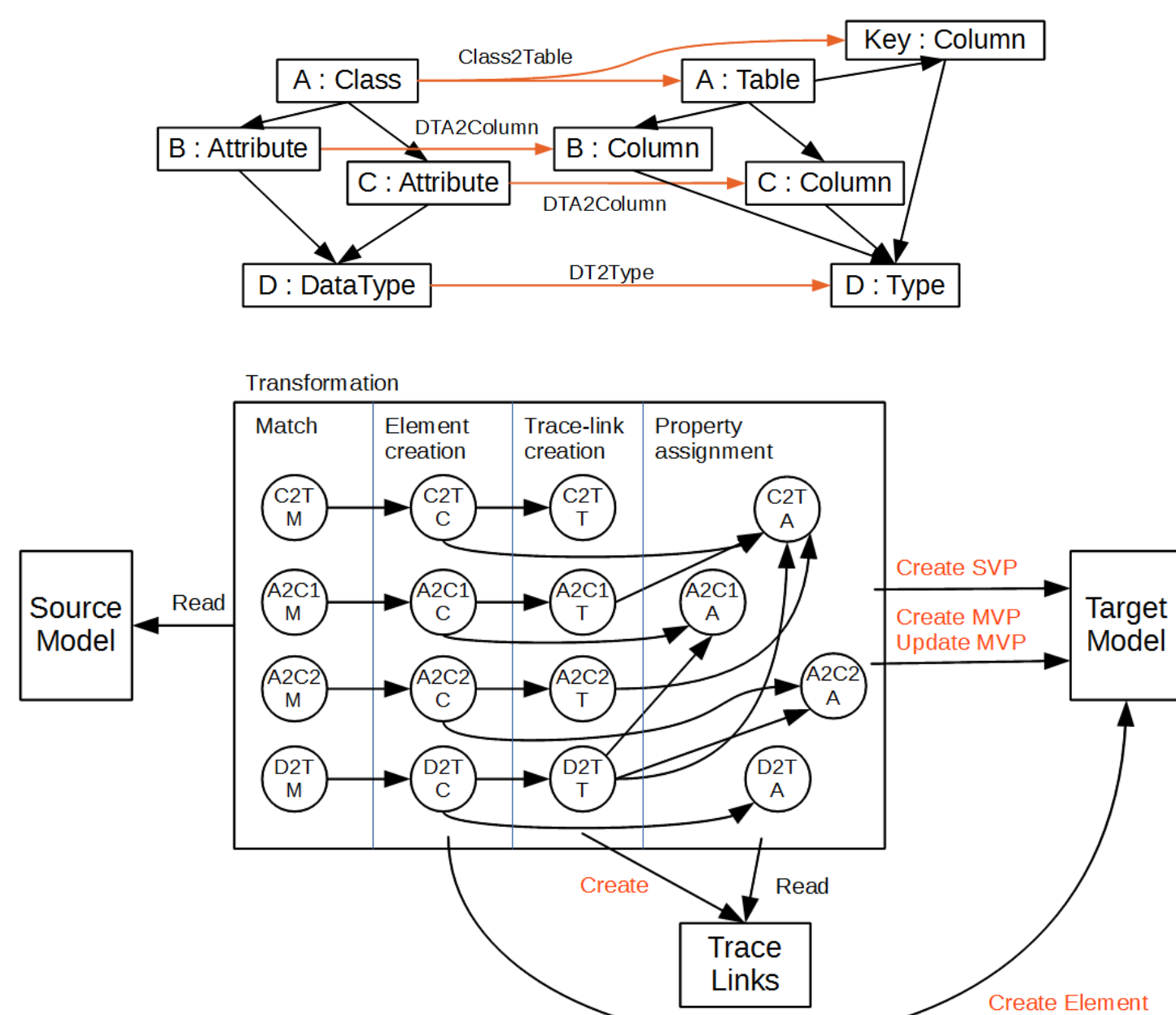



I- Model Driven Engineering

- Applying modeling to assist software development during the entire life-cycle
- Automating software development throughout code generation, validation, visualization etc.
- Large MDE community around Eclipse Modeling project
- Large set of tools in the Eclipse Modeling Project from industry and academics
- **The current generation of modeling technologies is stressed to its limits**
- **Examples from industry :**
 - ❑ Reverse-engineering systems with millions of LOC
 - ❑ Synchronizing views on building information models of several Gbs
- **Need of scalable MDE solutions for very large and complex systems**



II- Enabling scalability in Model Transformations (MT)

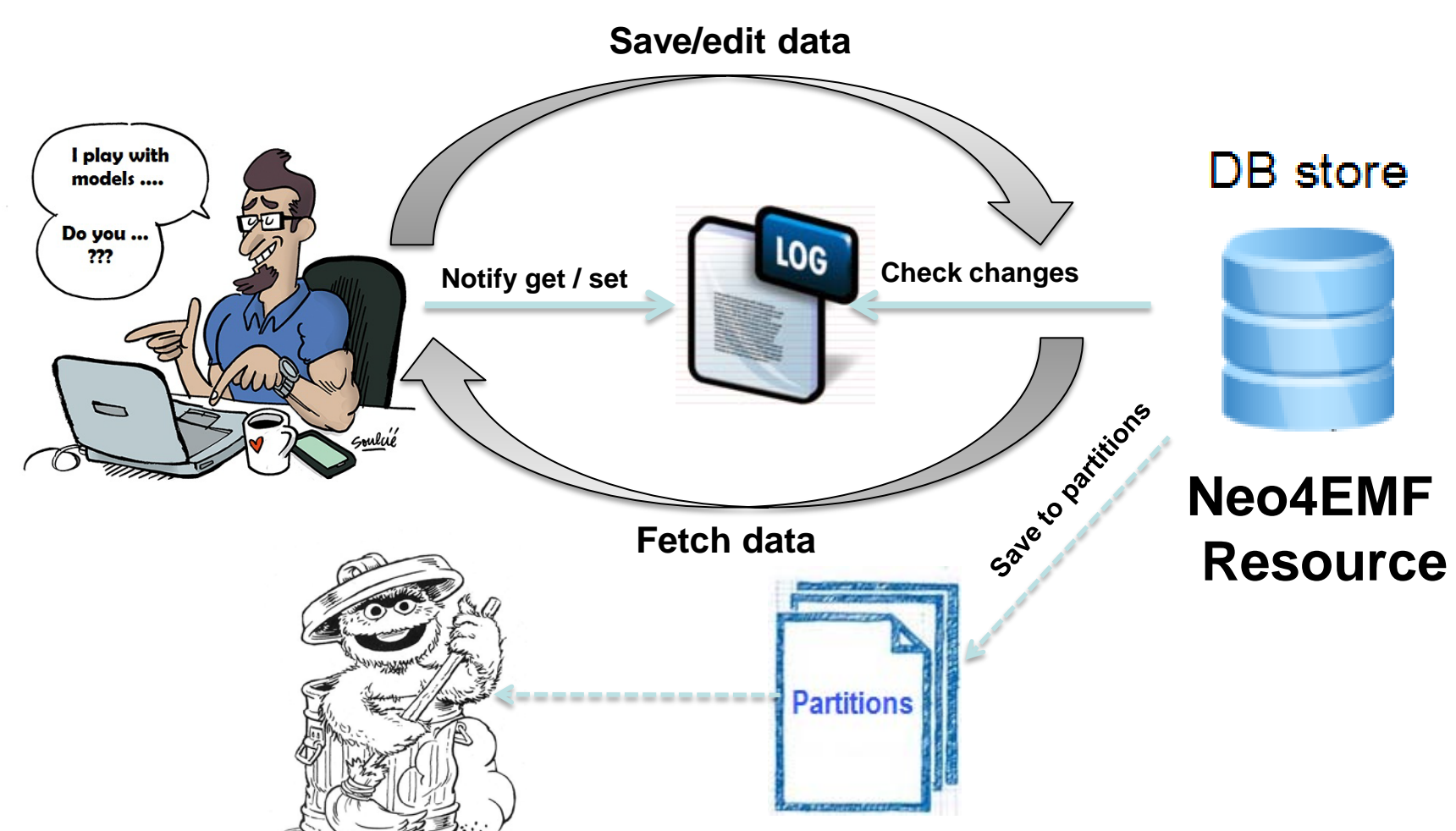


- The ATL transformation language:
 - ❑ is a model transformation language and toolkit developed in the AtlanMod research team
 - ❑ provides a parallel engine for faster transformations
 - ❑ enables change propagation and model synchronization using an incremental execution
 - ❑ reduces memory footprint and computation using lazy transformations
 - ❑ enables infinite transformations using lazy transformations

III- Enabling scalability in Model Persistence

Neo4EMF

- A model persistence framework (MPF) is a middleware that assists the storage of models
- The Neo4EMF model persistence framework:
 - ❑ provides a No-SQL backend using a graph database
 - ❑ enables loading large models using an on-demand loading mechanism
 - ❑ enables a lightweight first time loading by fetching objects not their data
 - ❑ involves a change (access) log to unload (save) models elements



Predicting personalized response to drug and environmental chemicals from genomic data with machine learning

1,2 Erwan Scornet, 1 Elsa Bernard, 1 Yunlong Jiao, 1 Veronique Stoven, 1 Thomas Walter, 1 Jean-Philippe Vert

¹Center For Computational Biology, Mines ParisTech ; INSERM U900 ; Institut Curie

²Paris VI

March 2014

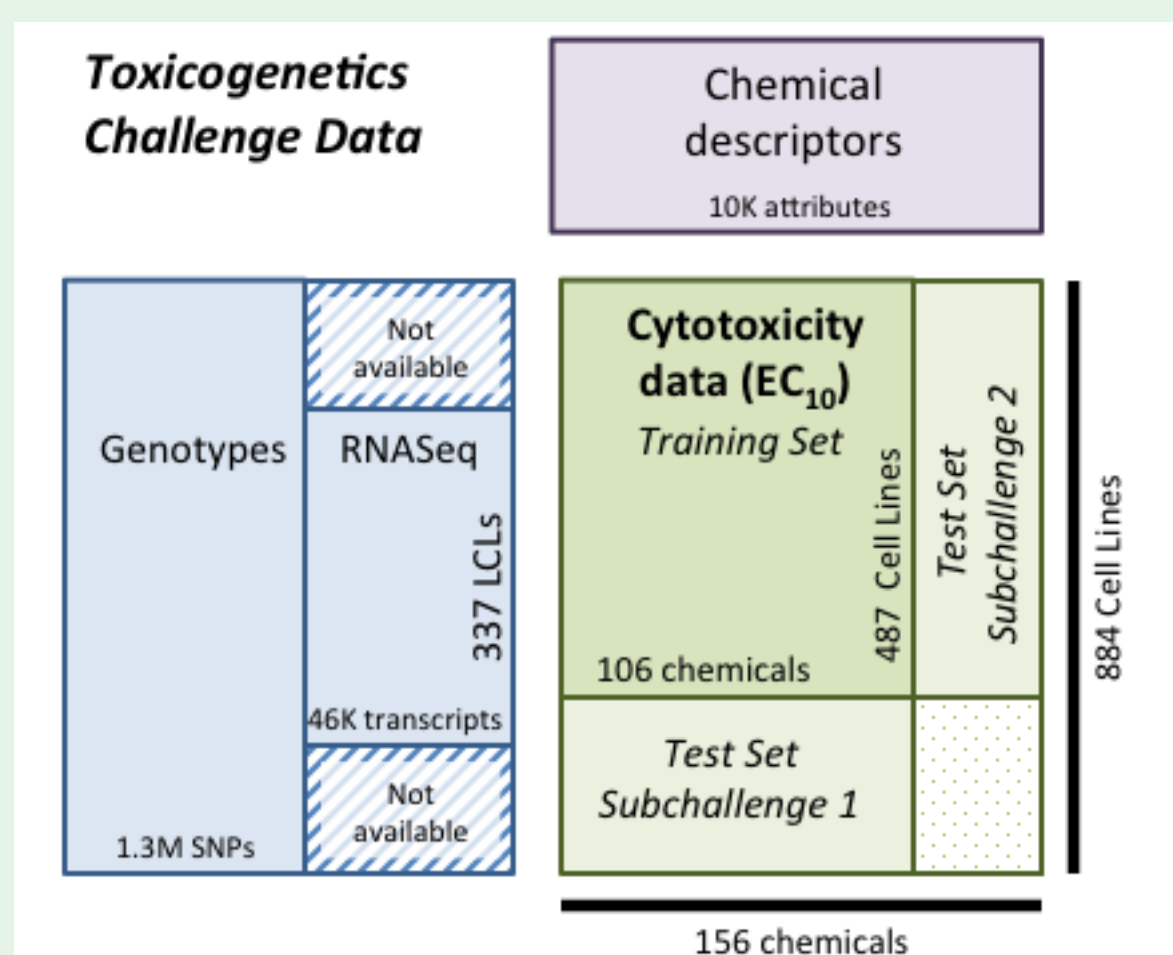


Abstract : The response to drugs and environmental chemicals changes with people. Variation in genotype can explain different reaction to drug : while some treatment may be effective for some people, it can be useless and even harmful for others. The same mechanisms are at stake regarding the reaction to environmental agents, such as allergens. Recent advances in high-throughput sequencing open the way to personalized treatment based on genotype data. In this work, we predict the toxicity level of 106 chemicals for each patient using both information on patients and on chemicals.

Challenge of Toxicogenetic

- Provide personalized information about chemical toxicity for each patient
- Use genetic information of each patient to predict the chemical toxicity
- Use chemical information (substructures, compounds,...) to learn toxicity across chemicals

Data description



Data were available thanks to the DREAM challenge 8.0 which gathers :

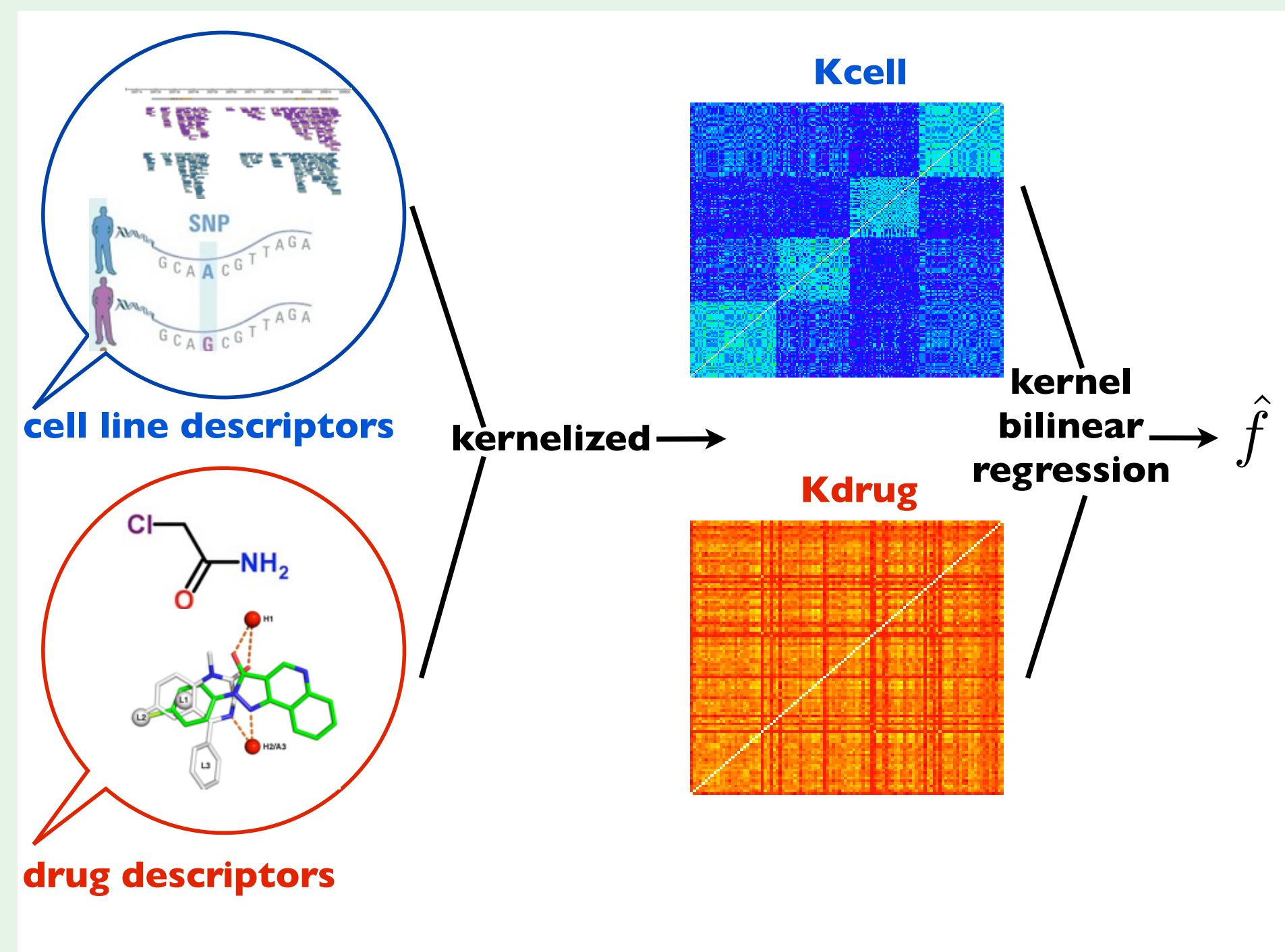
- Genotype data from the 1000 Genomes project (aiming at finding most genetics variants with frequency of 1%)
- RNA sequencing data from the Geuvadis Project (sequencing of some cell lines of the 1000 Genomes Project)
- Three covariates (sex, population and batch)

Toxicity → The drug concentration that reduces the ATP synthesis of 10%.

Methods

To predict the toxicity values, we aim at

- Creating a measure of similarity between chemicals
- Creating a measure of similarity between genotypes



Let \mathbf{x}_i be a vector of descriptors of cell line i and \mathbf{y}_j be a vector of descriptors of chemical j

- We model the toxicity t_{ij} of the chemical j on the cell line i by

$$t_{ij} = \underbrace{f(\mathbf{x}_i, \mathbf{y}_j)}_{\text{"bilinear"}} + \underbrace{b_j}_{\text{offset}} + \underbrace{\epsilon_{ij}}_{\text{noise}}$$

- We estimate \mathbf{f} and \mathbf{b} by penalized least-square regression

$$\min_{\mathbf{f} \in \mathcal{H}, \mathbf{b} \in \mathbb{R}^p} \sum_{i=1}^n \sum_{j=1}^p \{t_{ij} - f(\mathbf{x}_i, \mathbf{y}_j) - b_j\}^2 + \lambda \|\mathbf{f}\|_{\mathcal{H}}^2$$

To specify $\|\mathbf{f}\|_{\mathcal{H}}^2$ we just have to choose two kernels :

- $\mathbf{K}_{\text{cell}}(\mathbf{x}_i, \mathbf{x}_{i'})$ which measures the similarity between cell lines \mathbf{x}_i and $\mathbf{x}_{i'}$.
- $\mathbf{K}_{\text{drug}}(\mathbf{y}_m, \mathbf{y}_{m'})$ which measures the similarity between chemicals \mathbf{y}_m and $\mathbf{y}_{m'}$.

- Then, the solution of the previous optimization problem is given by

$$\hat{\mathbf{f}}(\mathbf{x}, \mathbf{y}) = \sum_{i=1}^n \sum_{j=1}^p \hat{\alpha}_{ij} \mathbf{K}_{\text{cell}}(\mathbf{x}, \mathbf{x}_i) \mathbf{K}_{\text{drug}}(\mathbf{y}, \mathbf{y}_j)$$

where $\hat{\alpha}_{ij}$ depend only on

- $\mathbf{K}_{\text{cell}} = (\mathbf{K}_{\text{cell}}(\mathbf{x}_i, \mathbf{x}_{i'}))_{i,i'=1,\dots,n}$
- $\mathbf{K}_{\text{drug}} = (\mathbf{K}_{\text{drug}}(\mathbf{y}_m, \mathbf{y}_{m'}))_{m,m'=1,\dots,p}$
- The toxicity matrix.

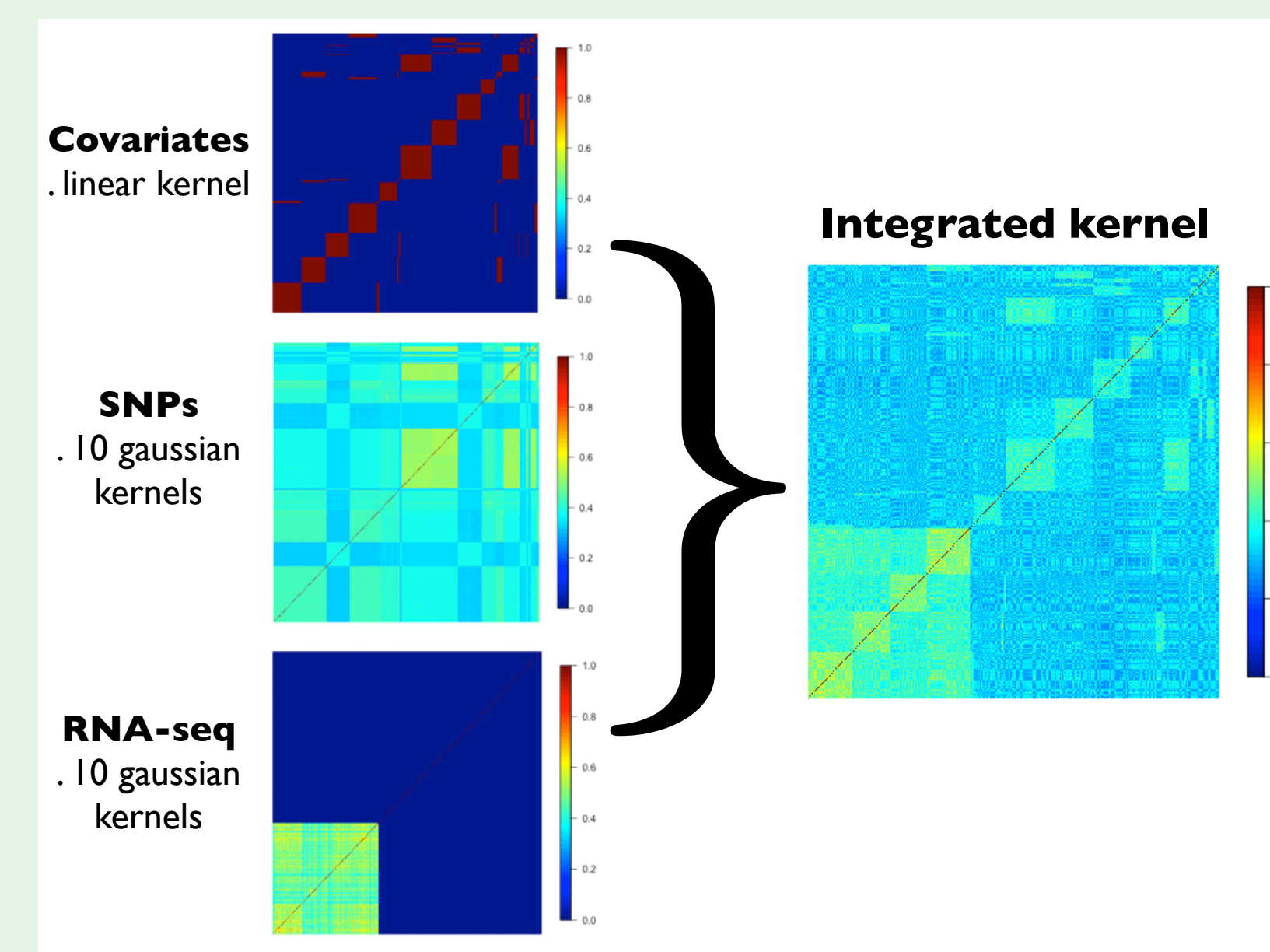
We also know that $\hat{\mathbf{b}}$ only depends on these three matrices.

References

- 1 B. Chem and al. *Journal of Chemical Information and Modeling* 49 :2044–2055, 2009.
- 2 Y. Yamanishi and al. *Journal of Chemical Information and Modeling* 51 :1183–1194, 2011.
- 3 P. Mahé and al. *Journal of Chemical Information and Modeling*, 45 :939–951, 2005.
- 4 T. Evgeniou and al. *Journal of Machine Learning Research* 6 :615–637, 2005.

What type of kernel do we choose ?

Kernel on cell lines :



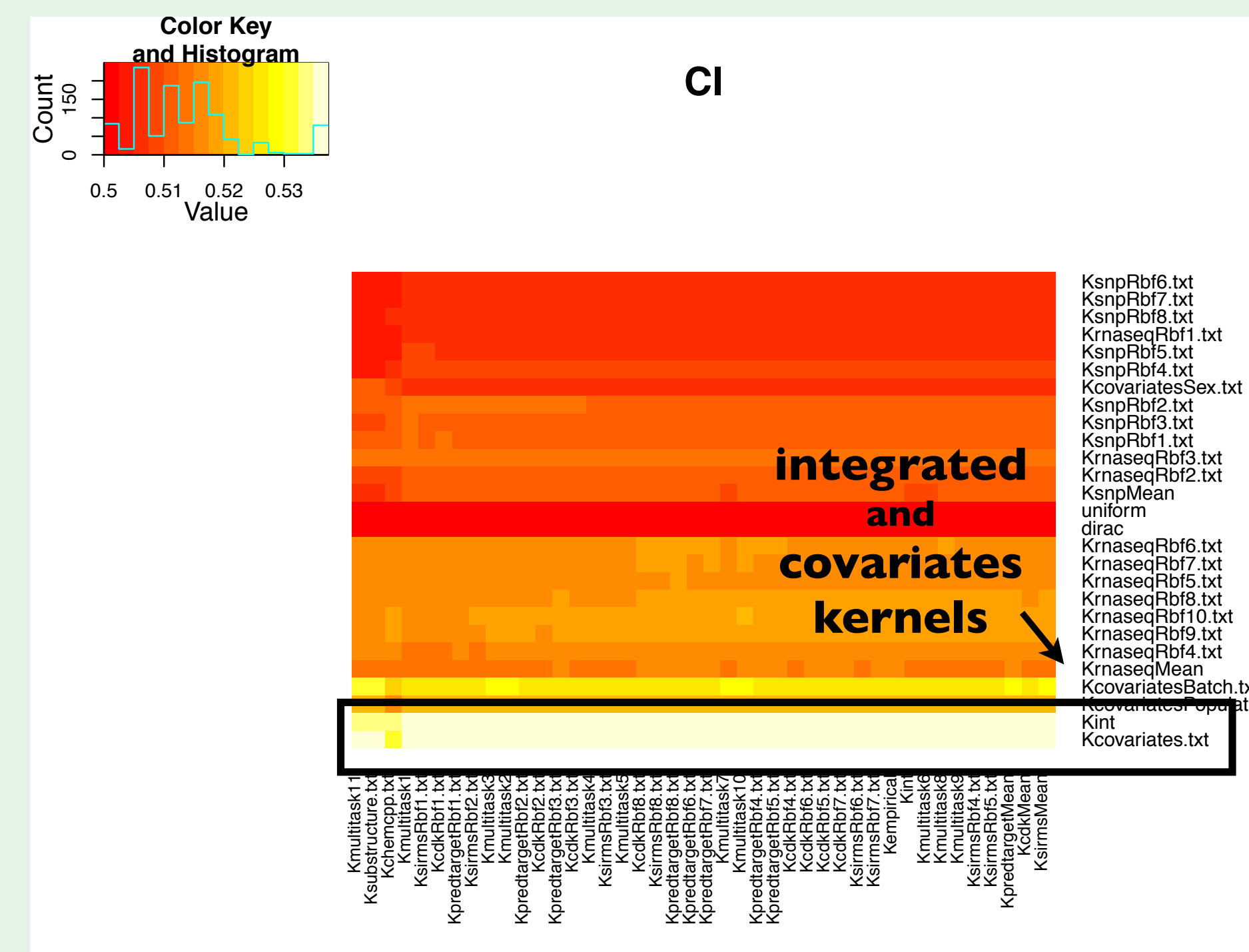
→ 29 kernels whose one that incorporate all informations (integrated kernel).

Kernel on chemicals :

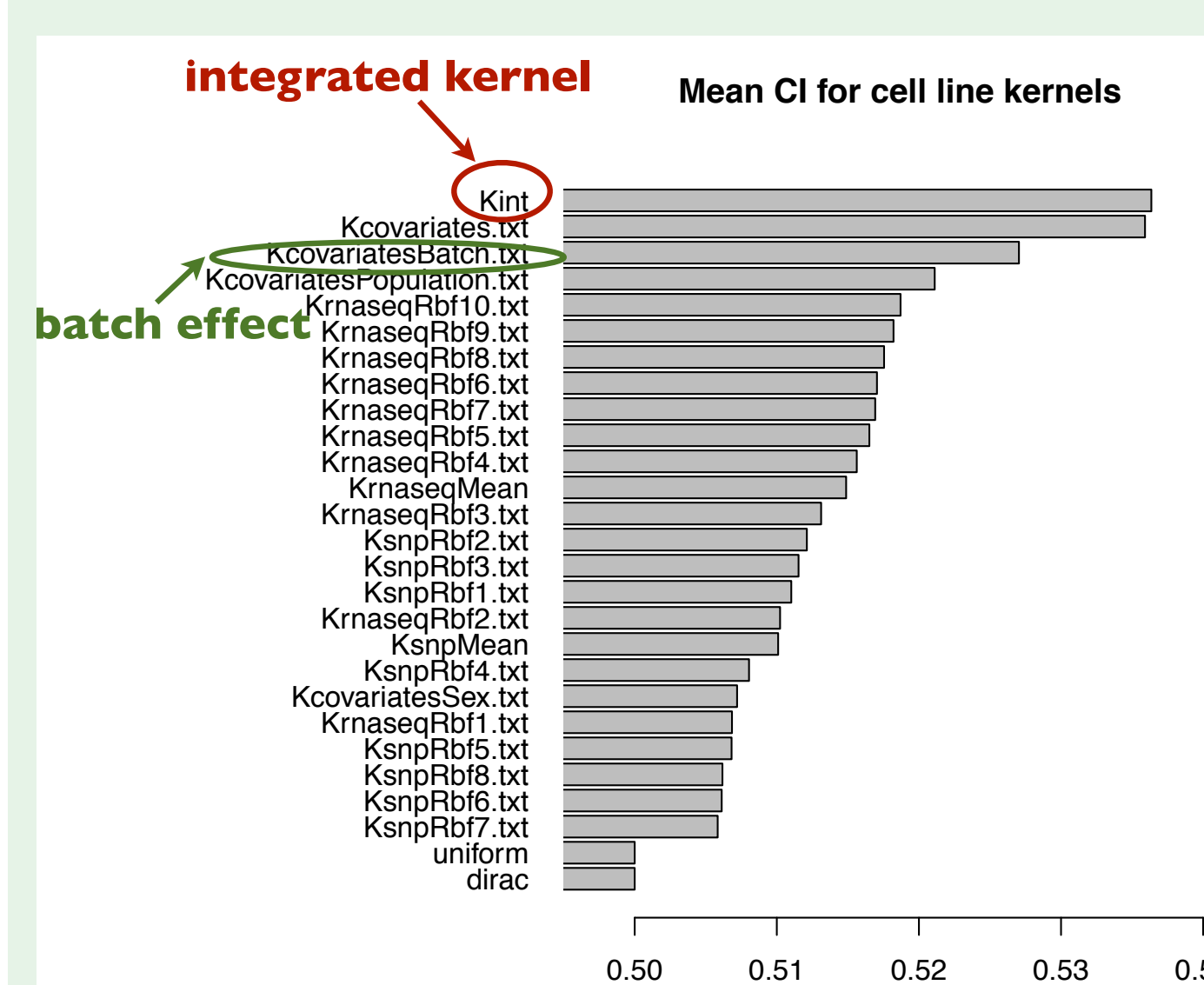
- Dirac kernel (build one model per chemical)
- Uniform kernel (build one single model)
- Kernel based on chemical features (using link to human protein, 2D substructures,...)
- Empirical kernel (evaluate the link between chemicals)
- An integrated kernel

Results

→ Evaluate the performance of each model by cross validation

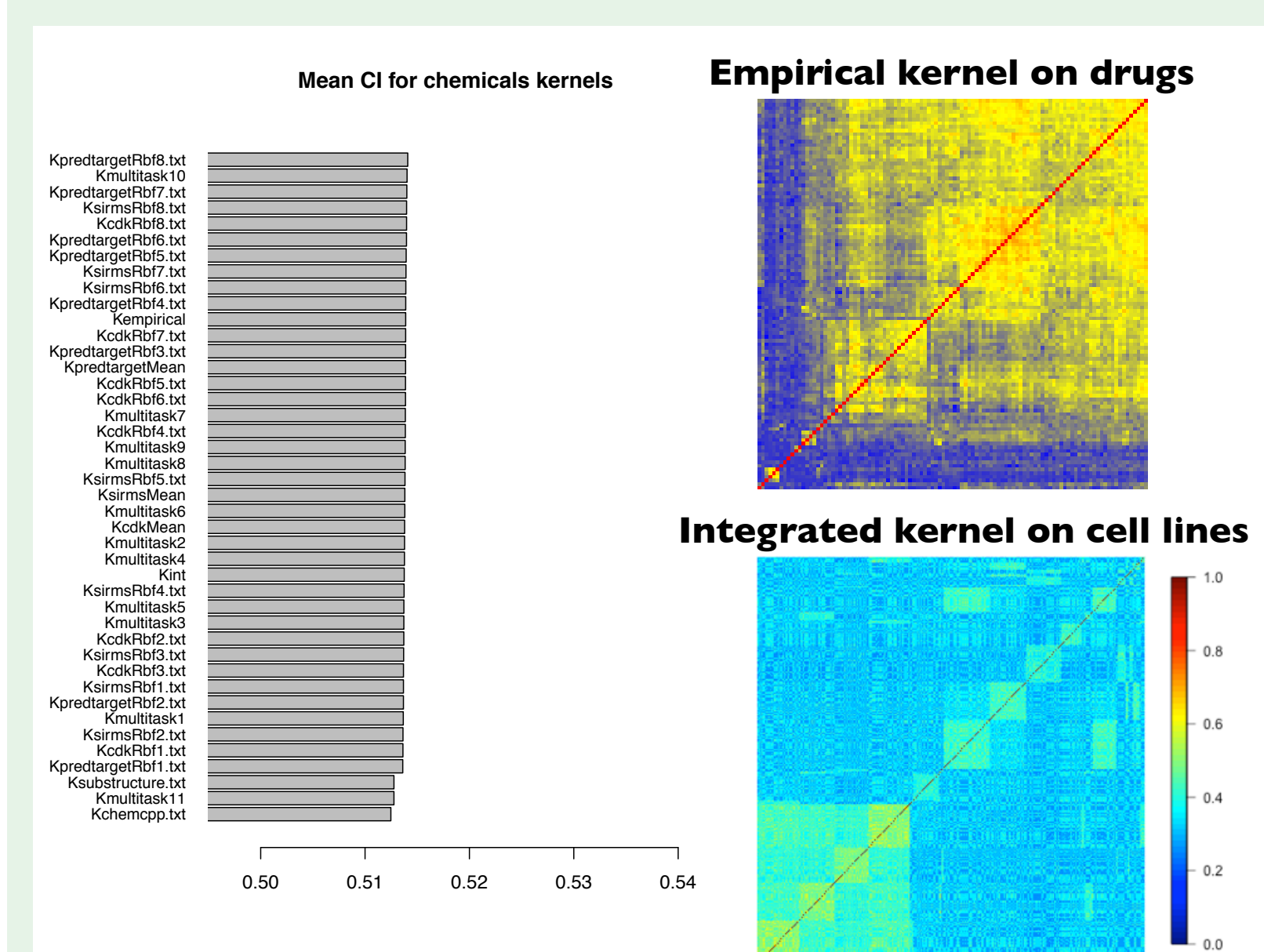


→ Integrated kernel and covariate kernel have the lowest concordance index.



- The integrated kernel has the highest prediction accuracy.
- The covariate kernel has a high prediction accuracy whereas it does not use any genomic information.
- This is due to a strong batch effect.

→ We choose the integrated kernel for cell lines.

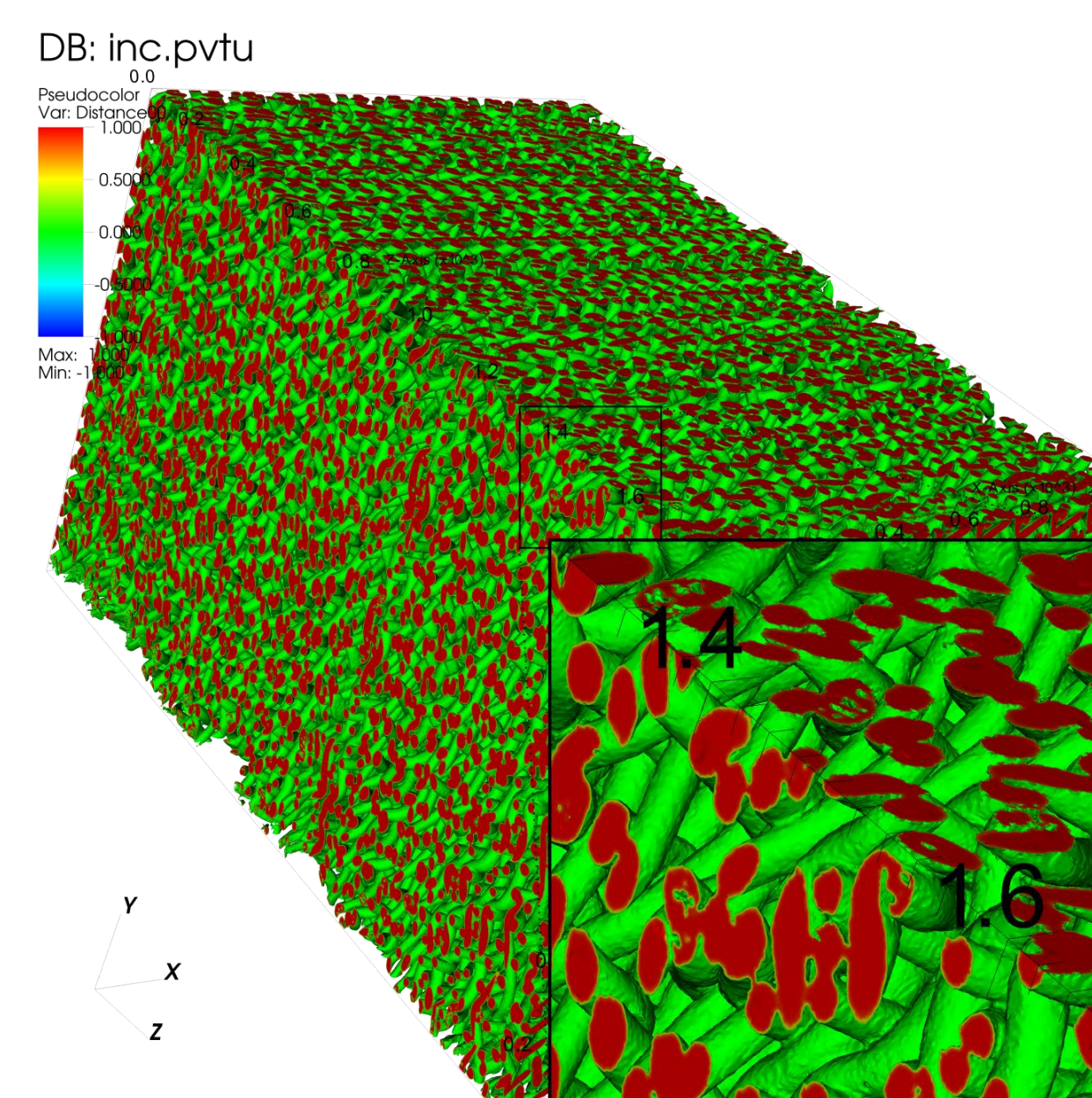


- All chemical kernels have the same accuracy : Dirac kernel performs as good as kernels using information on chemicals.

- Maybe the chemical descriptors do not fit this particular toxicogenetic problem.

→ We choose the empirical kernel for chemicals : since no kernel seems to be the best, we evaluate the link between chemicals instead of using some predetermined kernel.

→ This model ranked second out of 100 models submitted, in the DREAM challenge 8.0.



1200x1200x1791
780 millions de nœuds
4000 cœurs de calcul

STRATEGIE GENERALE

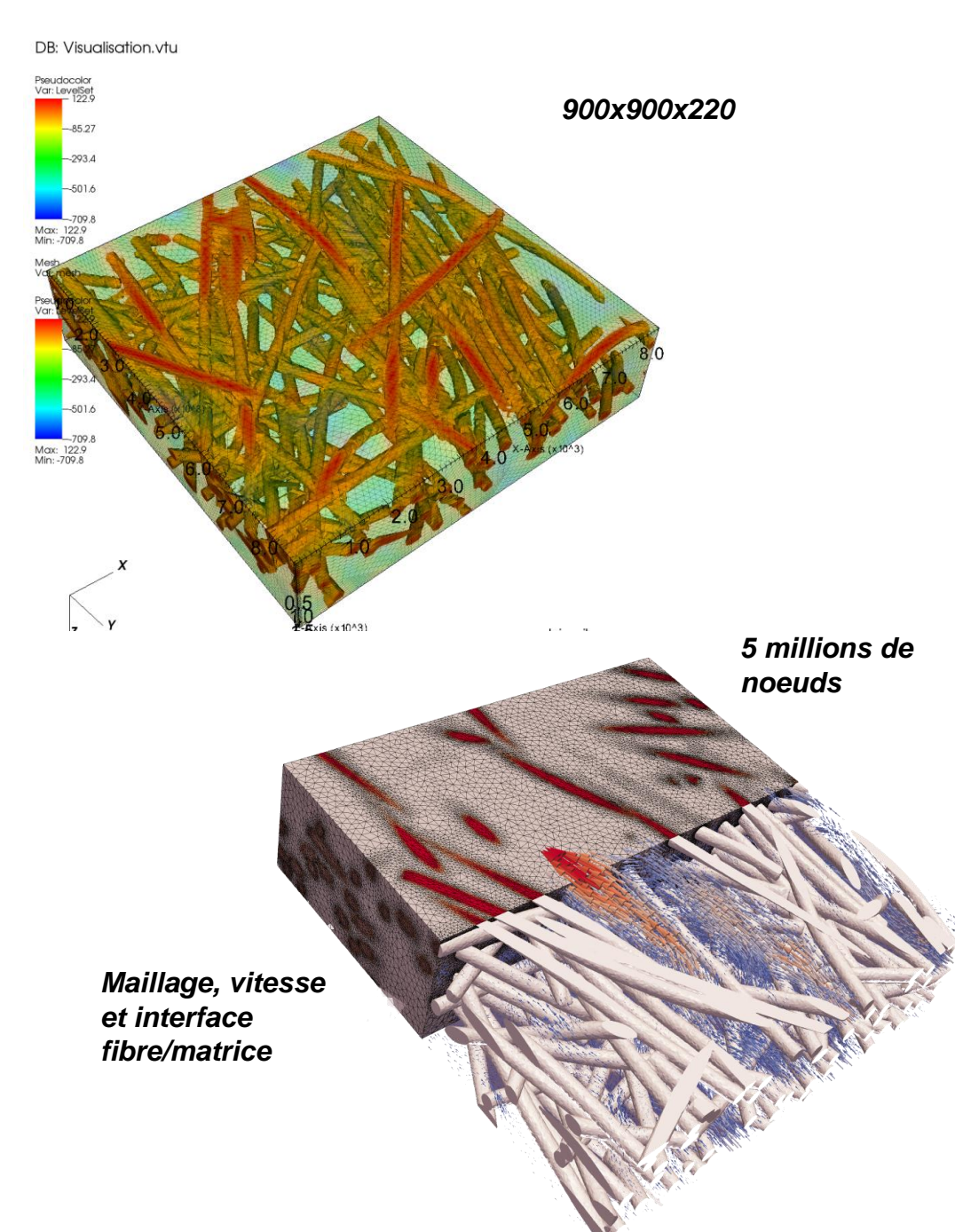
Objectifs

- Construire des maillages de façon pertinente et efficace à partir d'imagerie 3D
- Compression de l'image, sans perte d'information, par manipulation au format maillage
- Utilisation de ces maillages dans différents domaines d'application

Résultats atteints / prototypes / démonstrateurs

- « Maillage d'images »: génération avec un maillage topologique et adaptation de maillage par minimisation de l'erreur d'interpolation de l'image sur le maillage, associée à une technique de réinitialisation de la valeur du voxel/pixel
- Au bout d'une année: génération de microstructures virtuelles, de géométries élancées en sous-sol, d'environnements urbains; calculs éléments finis et éléments frontière sur certaines applications
- Plateformes logicielles: CimLib, Morph-M, Neper, Zebulon, ...

Parties prenantes



MATERIAUX COMPOSITES

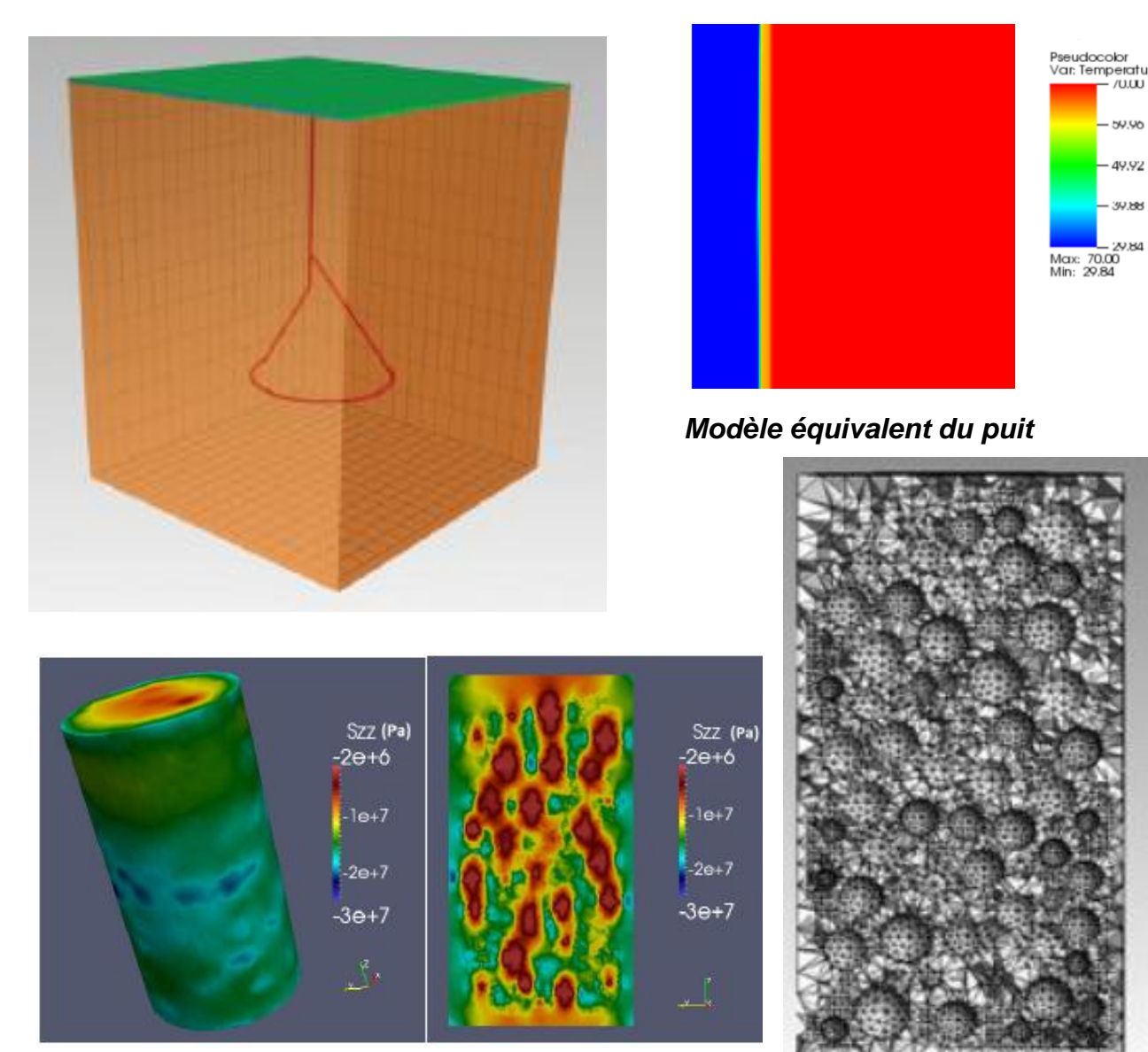
Perméabilité d'un milieu fibreux

- Arrangement 3D irrégulier obtenu par imagerie 3D [Orgéas et al, 3S-R], caractéristiques renfort: $R=0.1\text{mm}$, $L=10\text{mm}$, $\phi=0.83$
- Images acquises par microtomographie-X (900x900x220 voxels), taille du voxel = $10\ \mu\text{m}^3$
- Génération du maillage éléments finis (~5 millions de nœuds): interpolation directe de l'image et adaptation, sur 96 cœurs
- Calcul d'écoulement sur le système fibre- matrice (sur 96 cœurs) et homogénéisation pour obtention de la perméabilité

MICROSTRUCTURES METALLIQUES

Déformation d'un polycristal

- Image obtenue par tomographie de contraste à diffraction (DCT)
- Adaptation d'un algorithme de génération de cellules de Voronoï aux grains de l'image et maillage associé avec Neper
- Calculs de déformation sur l'échantillon avec Zebulon



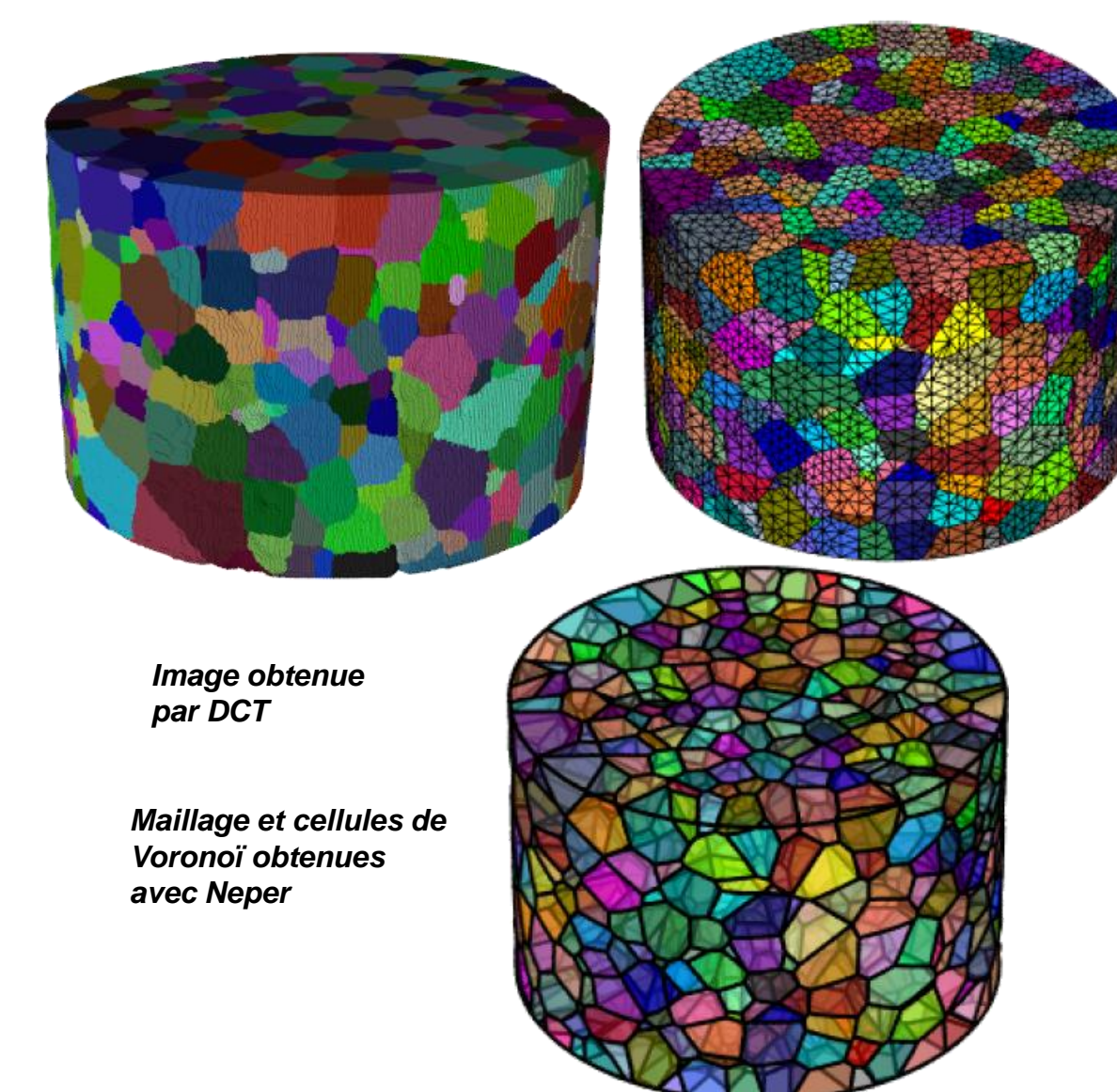
GEOSCIENCES

Modélisation d'un échangeur géothermique

- Ecoulement d'eau dans un tube avec gainage (et couplage avec la température), validation CimLib/Fluent en 2D pour un modèle simple

Modélisation de la méso-structure des bétons légers

- Compression d'un béton avec un calcul à l'échelle de la méso-structure pour obtenir le module de Young et proposer un mécanisme pour la rupture des bétons légers



Auteurs

Luisa Silva
Rima Ghazal
Daniela Craciun
Jia-Xin Zhao
Min Quan Thai
Sébastien Nadler
et al.

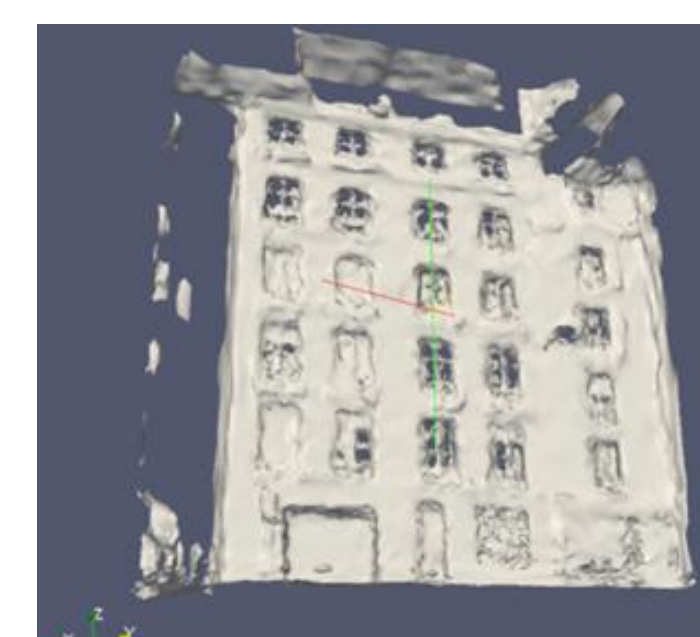
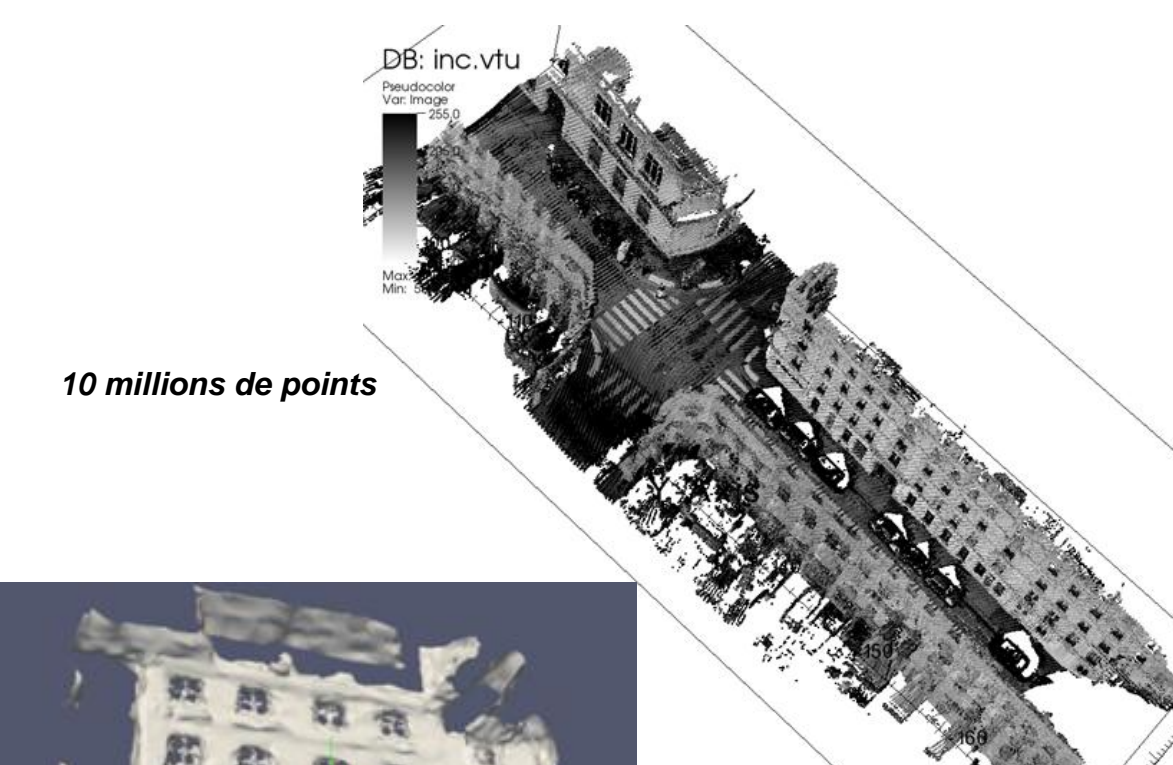
Partenaires



MODELISATION DES ENVIRONNEMENTS URBAINS

Génération de maillages 3D à partir de nuages de points

- Rue de Paris, capture obtenue avec un Velodyne, 10 millions de points
- Compression de l'information par utilisation d'un maillage 3D surfacique, mais aussi 3D volumique (immersion dans une géométrie volumique)
- Applications: calculs sur ces environnements



MoGDIW, an integrated workflow for cell motility genes discovery in high-throughput time-lapse screening data



^{1,2,3}Alice Schoenauer Sebag, ²Céline Raulet-Tomkiewicz, ²Robert Barouki, ¹Jean-Philippe Vert, ¹Thomas Walter

¹Center For Computational Biology, Mines ParisTech ; INSERM U900 ; Institut Curie

²INSERM U747 ; Paris V

³Agro ParisTech

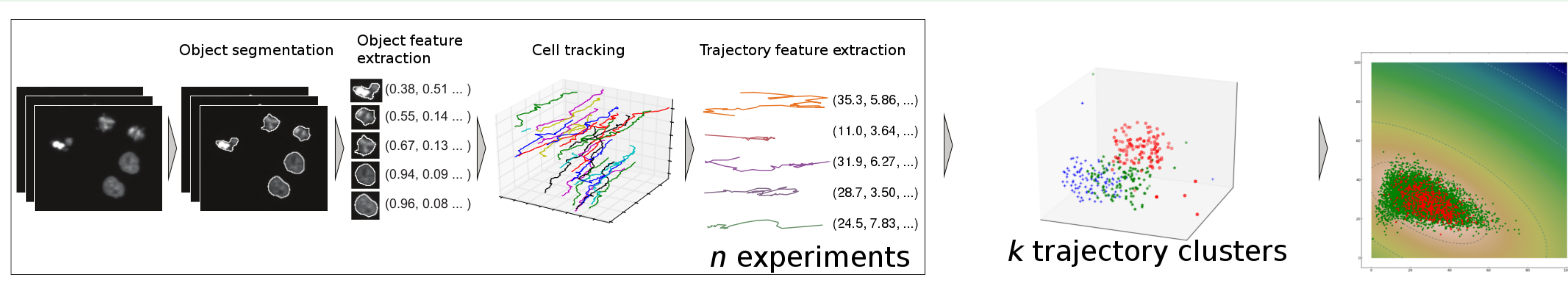
December 2013



Abstract : Cellular migration is a fundamental biological process. Progress in the fields of gene silencing and high-throughput (HT) microscopy has only recently made its study possible on a large scale. However, all existing HT migration screens measure motility at the level of cell population. Here, we present MoGDIW, a generic integrated workflow which addresses cell motility genes discovery in HT time-lapse screening data at single cell level. It is composed of cell tracking, cell trajectory mapping to an original feature space, migration pattern identification, and discriminant characterization of each experiment in terms of migratory behaviours. In comparison with an existing migration screen, MoGDIW application to a genome-wide time-lapse screen shows little overlap. However, its results are enriched in migration and adhesion-related (MAR) genes, and could be visually confirmed.

Motility Genes Discovery Integrated Workflow

Aim : quantitatively assess and compare single cell migration under different chemical perturbations in time-lapse microscopy data



Cell tracking

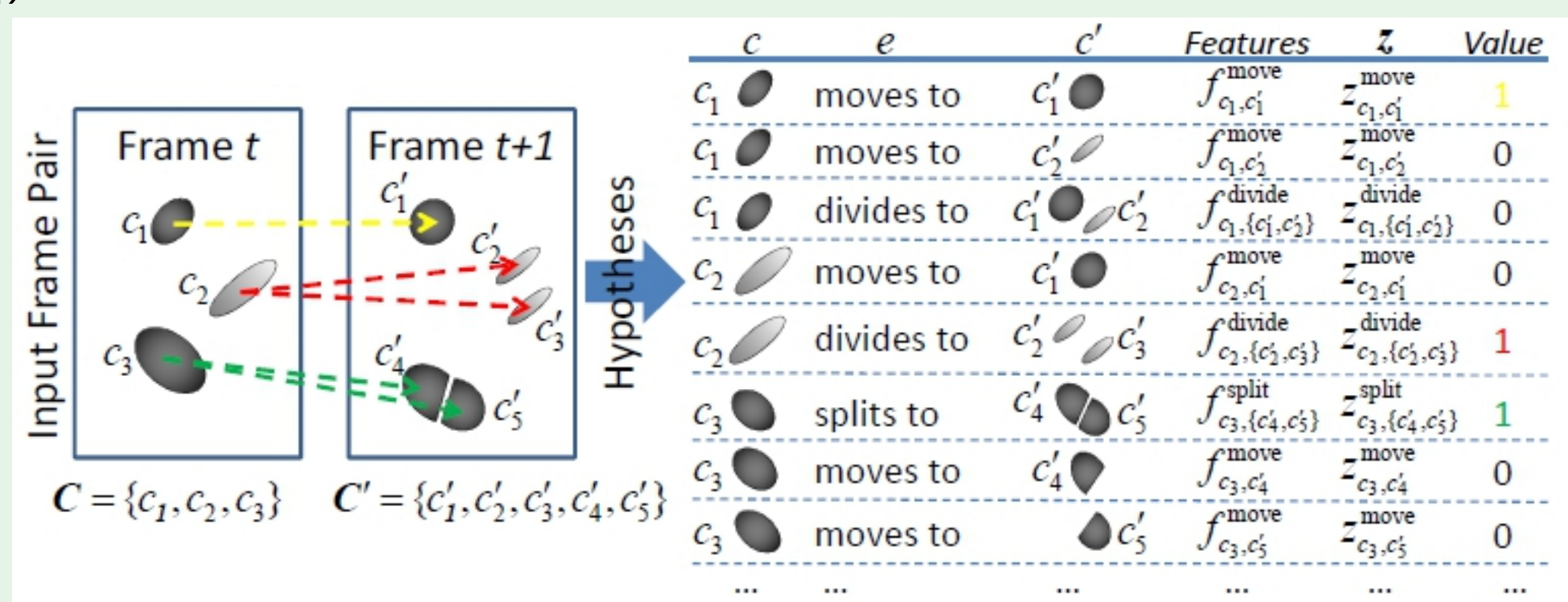
Challenges :

- High population density and high phenotypic variability
- Frequent absence of overlaps in consecutive frames (time interval : 30')
- Appearing, disappearing, merging and splitting objects
- Minimum use of prior knowledge
- Usable for non-experts

⇒ Strategy : use of Machine Learning

Structured learning for cell tracking

- Idea (from [3]) :



- Prior knowledge : possible events **e** which can occur to an object between two frames
{move, appear, disappear, split in 2 or 3, merge at 2 or 3}

- Data : consecutive frames with already segmented objects

- Target : learn the model on annotated data and use it.

- Model :

$$\hat{z}(t) = \arg \max_z L(x(t), z; w) = \sum_{\substack{e \\ \text{Obj}_{j,t} \\ \text{Obj}_{j,t+1}}} \langle w^e, f_{i,j}^e \rangle z_{i,j}^e$$

$$\text{st } \forall i \in \{0, \dots, N(t)\} \sum_{\substack{e \\ \text{Obj}_{j,t+1}}} z_{i,j}^e = 1$$

$$\text{and } \forall j \in \{0, \dots, N(t+1)\} \sum_{\substack{e \\ \text{Obj}_{j,t}}} z_{i,j}^e = 1$$

- Learning **w** : Support Vector Machine (SVM, algorithm : bundle method [6])

Results

- Training set : ~ 32 000 links with 0.5% appear ; 0.5% disappear ; 1% merge ; 2% split

Algorithm (software)	Mean recall	Mean precision	Ref
Constrained nearest neighbour (Cell Cognition)	72.7%	62.8%	[1]
Linear assignment problem (Cell Profiler)	78.3%	73.0%	[2]
Structured learning	91.1%	91.5%	[3]

Providing a graphical user-interface for annotating videos

- Integration in Cell Cognition [1], an open source software platform for the analysis of live cell imaging data, with the IMBA, Vienna

- Extensions with regard to cell track annotation :

- Generation of cell tracks rough estimation using a Nearest Neighbor tracker
- GUI extension to support manual correction of the Nearest Neighbor trajectories

References

- Held et al., Nature Methods, 7(9) :747-54, 2010.
- Jaqaman, K. et al, Nature Methods, 5(2008) :695-702.
- Lou et Hamprecht, NIPS, 2011.
- Neumann et al, Nature, 464 :721-727, 2010.
- Simpson et al, Nat. Cell Biol., 10 :1027-1038, 2008.
- Tsochantaridis et al, JMLR, 6 :1453-1484, 2005.

Experiment characterization in terms of migratory patterns

Trajectory mapping to an original feature space :

TABLE: Feature types and examples of corresponding features

Goal	Examples of feature
Track characterization	Diffusion coefficient, persistence, track evenness
Static quantification	Convex hull area, largest move, total path length
Dynamic quantification	Mean acceleration, mean instantaneous speed

Experiment characterization :

- Clustering of all trajectories across experiments
- Characterization of an experiment by the cluster histogram of its trajectories
- Pearson's χ^2 test for testing significant deviation from control histogram

Analysis of Mitocheck screen

- Data : MitoCheck data set [4], ~ 200,000 videos of HeLa cells, produced by high-throughput live cell imaging, following selective down-regulation of all protein coding genes, one by one, by RNA interference (RNAi)

- Gene subset : 1,081 genes previously selected by [5] consisting of

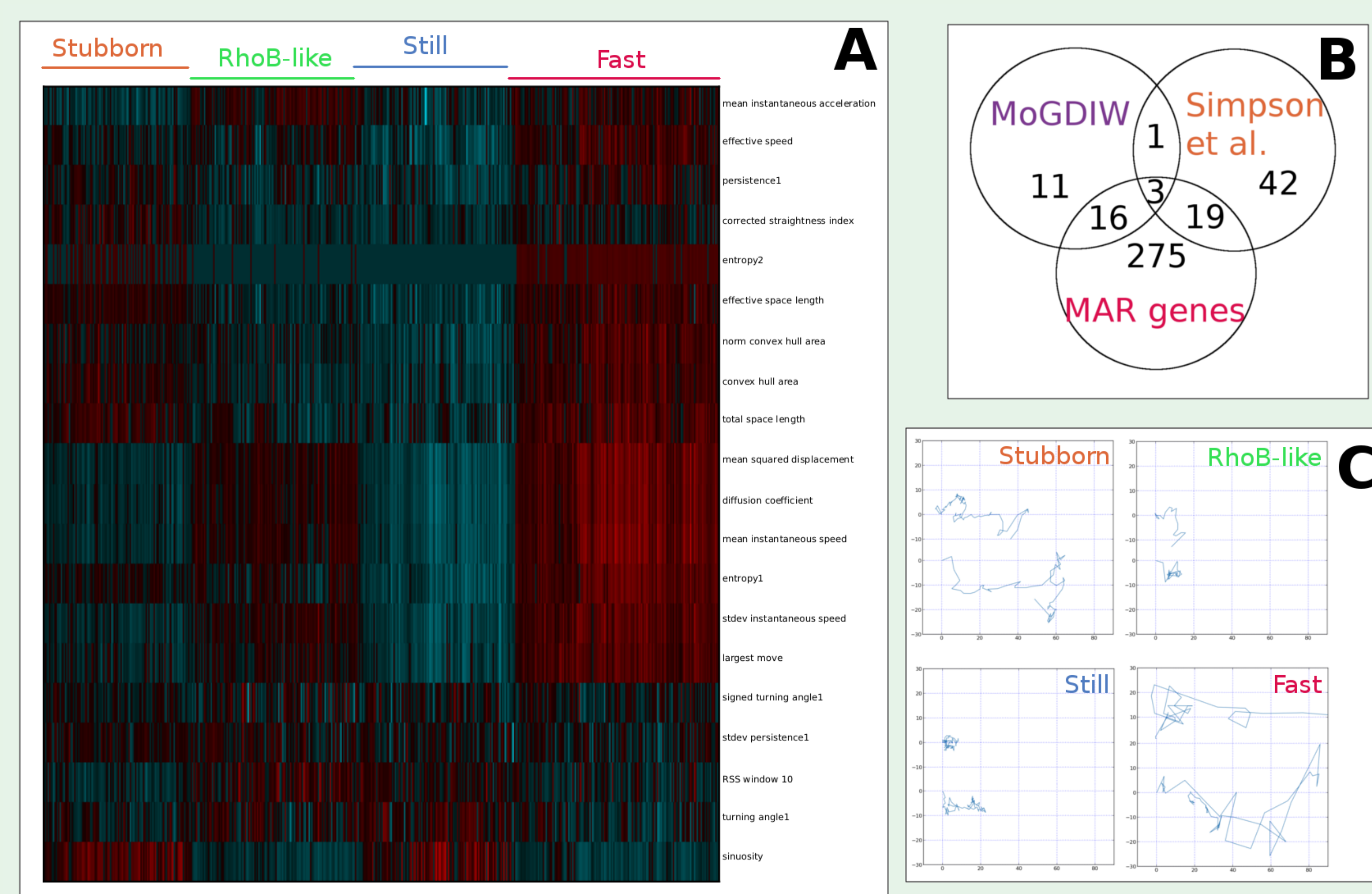
- phosphatases
- kinases
- migration and adhesion-related genes (MAR genes, a priori selected by the Geiger Laboratory, Weizmann Institute, Israël)

- Method : MoGDIW with K-means, k=4

- Results in comparison with [5] :

- Small overlap of selected gene lists
- No enrichment in known migration genes in either case
- Enrichment of MoGDIW high-confidence list in MAR genes (61% vs 34% for [5])

FIGURE: A. Identification of MoGDIW cluster characteristics using single linkage hierarchical clustering ; B. Venn diagram comparing MoGDIW's and [5]'s high-confidence genes ; C. Trajectory examples



Perspectives

- Application to the whole Mitocheck dataset and biological validation of hit genes
- Application to newly generated Environmental Toxicology data to assess the consequences of chemical exposure on single cell migration

Adaptation de maillage anisotrope, calcul parallèle et capture d'interface Applications au matériaux et à la mise en forme, mais pas seulement...

PARTENAIRES

Entre autres...



PUBLICATIONS

Quelques unes...
 Digonnet, *App Math Model* (2000)
 Bruchon, *Int Journ Num Meth Eng* (2009)
 Hachem, *Journ of Comp Physics* (2010)
 Coupez, *Journ of Comp Physics* (2011)
 Ville, *Int Journ Num Meth Fluids* (2011)
 Silva, *Int Journ Mat Form* (2012)
 Bernacki, *Scripta Mat* (2011)
 Tillier, *Int Journ Or Maxi Surg* (2012)
 Carozzani, *Met Mat Trans* (2013)

REMERCIEMENTS



CONTACT

luisa.silva@mines-paristech.fr
 www.cemef.mines-paristech.fr
 CEMEF MINES ParisTech
 CS 10207
 06904 Sophia Antipolis - France

CimLib

C++ et calcul scientifique

Caractéristiques phare

- une base éléments finis
- haut degré de parallélisme
- adaptation de maillage et du temps anisotrope
- méthodes de stabilisation performantes
- calcul de surfaces libres ou interfaces par des méthodes robustes
- validation automatique et mise à disposition sur un serveur de partage de projet
- utilisée pour le développement de logiciels industriels: Rem3D, Thost, Ximex, Forge, Thercast, Transweld

Adaptation anisotrope

Maillage et temps

Automatique, en parallèle et en dynamique

2h, 23 millions de nœuds, 10 itérations, 256 cœurs

Fonction:
 $f(x) = g + g(x-0)$
 $*g + g(x-1)$
 $g(x) = \tanh\left(\text{Erfm}\left(\frac{x-0.5}{\sigma}\right)\right)$

Simulation de l'injection multi-fluides et adaptation du maillage et du temps

Massivement parallèle

Supercalculateurs et supercalculs

Partitionnement dynamique, solveur multigrilles

Génération, adaptation de maillage			Résolution d'un très grand système linéaire			
dim	# éléments (milliards)	# cœurs	dim	# inconnues (milliards)	temps calcul	# cœurs
2D	66,8	100 000	2D	100	319 s	262 144
3D	81,4	100 000	3D	55	447 s	65 536

Méthodes numériques avancées

- ▶ Génération et adaptation de maillages anisotropes
 - Adaptation anisotrope en dynamique et en parallèle (3D) basée sur une carte de métrique et sur des estimateurs d'erreurs anisotropes 3D sur le gradient des fonctions ou sur la distribution de la longueur des arêtes, avec contrôle du nombre d'éléments
- ▶ Calculs massivement parallèles
 - Haute performance des simulations grâce à l'exploitation efficace de la parallélisation de la gestion du maillage et de la résolution des systèmes linéaires
 - Repartitionnement dynamique, solveurs itératifs parallèles et multigrilles
 - Benchmarking et applications dans des supercalculateurs du Tiers1 (Genci) et Tiers0 (Prace)
- ▶ Écoulements multiphasiques, interactions fluide-structure et interfaces mobiles
 - Approche monolithique et méthode des volumes immergés
 - Méthodes éléments finis stabilisées pour les écoulements du très petit (très visqueux) au très haut (peu visqueux) nombre de Reynolds, avec couplages thermiques et cinétiques
 - Interactions thermomécaniques entre les phases, changement de phase
 - Evolution des interfaces par des approches type level-set

Aérodynamique et aérothermie

Écoulements turbulents, transferts thermiques

Interactions avec les structures, fixes ou mobiles

Écoulement d'air Haut Re (10^8)

Drone léger à aile bouclée et dirigeable à 7000 m

Simulation et le vivant

Biomécanique

Tissus mous ou structures osseuses

Calcul mécanique et lois de comportement pertinentes

Reconstruction des géométries

Écoulement multiphasique

Interfaces implicites évolutives

Approche monolithique et méthode level-set

Evolutions complexes: chute d'une goutte de lait

Injection assistée-eau: transport multi level-set

Physique numérique

Changements de phase

Germination et croissance

Ebullition
 Croissance et ascension de bulles de vapeur

Recristallisation
 Croissance de grains

Solidification
 Croissance dendritique et anisotropie de l'interface

Simulation et le réel

Imagerie et Données Massives

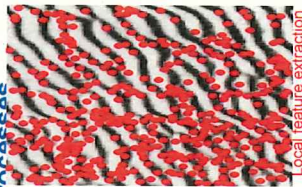
3D volumique ou surfacique

Microstructures numériques, biomédicale et tomographie-X

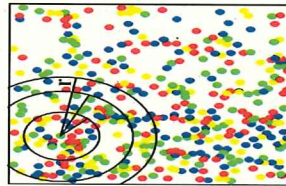
Environnements urbains et nuages de points

Abstract: In this paper, invariant texture characterization and recognition are addressed from the characterization of the spatial distribution of image. Visual keypoint sets in visual textures are here regarded as realizations of spatial point processes. We show that empirical second-order statistics considered in [4] relate to a non-parametric form of a log-Gaussian Cox model and investigate the relevance of parametric Cox models for texture recognition issues. Reported results validate the proposed descriptor compared to state-of-the-art approaches⁽¹⁾ with three datasets: UIUC, KTH-Tips, Brodatz.

Descriptive statistics of multivariate point processes



Local feature extraction



Codebook construction of keypoints⁽²⁾

A spatial point process S is defined as a locally finite random subset of a given bounded region B in \mathbb{R}^2 . A realization of such a process is a spatial point pattern $\{s = \{s_1, \dots, s_n\}$ of n points in B . Given realizations of a point process, the moments of random variable are relevant descriptive statistics. The p^{th} -order moment is defined by:

$$\mu^{(p)}(B_1 \times \dots \times B_n) = E\{N(B_1) \cdot \dots \cdot N(B_n)\}$$

For a marked spatial point process $\{(s_i, m_i)\}$ in given bounded regions B_i , where m_i is a mark associated to point s_i [3], the second-order moment is given by:

$$\alpha^2(B_1 \times B_2) = E \sum_{s_i \in \mathcal{S}(B_1 \neq B_2)} \sum_{s_j \in \mathcal{S}(B_2)} \alpha_i^2(s_i) \alpha_j^2(s_j)$$

Ripley's K function considers circular analyzing regions and resorts to the mean numbers of points of type j in a region of radius r centered at the points of type i :

$$\begin{aligned} K_{ij}(r) &= (\lambda_i \lambda_j)^{-1} \alpha_{ij}^2(r) \\ &= (\lambda_i \lambda_j)^{-1} E \sum_{h \neq i} \sum_{h \neq j} \delta_i^h(m_h) \delta_j^h(m_h) \mathbb{1}(\|s_h - s_i\| \leq r) \end{aligned}$$

Log-Gaussian Cox

Cox processes $\{X_i\}$ with **intensity functions** $\{\lambda_i\}$ are point processes such that $X_i | Z_i$ is a Poisson process, where $Z_i = \exp(Y_i)$. $\{Y_i\}$ is a multivariate Gaussian field on S characterized with mean function $\mu = EY(s)$ and covariance function $c_{ij}(r) = \text{Cov}(Y_i(s_i), Y_j(s_j, r))$ [4].

- Intensity function: $\lambda = \exp(\mu + \sigma^2/2)$
- Pair correlation function: $g_{ij}(r) = \exp(\rho_{ij}(r))$
- Ripley's K function: $K_{ij}(R) = 2\pi \int_0^R g(r) r dr$

The estimation of the pair correlation function is given by:

$$g_{ij}(r) = (2\pi \lambda_i \lambda_j)^{-1} \sum_{h \neq i} \sum_{h \neq j} \delta_i^h(m_h) \delta_j^h(m_h) \mathbb{1}(\|s_h - s_i\| = r)$$

Given a parameterization $L(\beta, r)$, namely Exponential, Hyperbolic or Cardinal sine, model parameters are estimated from the minimization of the following criterion:

$$\int_0^R \left\{ \sigma_{ij} L(\beta_{ij}, r) - c_{ij}(r) \right\}^2 dr$$

The proposed descriptor is formed by: $(\lambda_{ij}, \sigma_{ij})$



Edge-effect corrections:

Feature dimension reduction: A codebook of keypoint pairs $u = M(s_i, s_j)$ from two categorized keypoint s_i, s_j is considered, such that:

$$g_{ij}(r) = (2\pi \lambda_i)^{-1} \sum_{h \neq i} \sum_{h \neq j} \delta_i^h(m_h) \delta_j^h(m_h) \mathbb{1}(\|s_h - s_i\| = r) \delta_{ij}^h$$

Scaling effects: The actual radius of image is estimated by a reference radius r_{ref} and scale factor ϕ_{ref} (the rate of average point densities per surface unit).

$$r_i = \frac{\phi_{ref}}{\phi_i} r_{ref}$$

Application to texture

Parameter setting: recognition

- Categories of visual keypoints $k = \{60, 120, 150\}$, pairs of keypoints $k^* = 60$.
- $r_{ref} = \min(w, h)/2 \ln(x)$ where $w(dth), h(ght)$ the size of image, $x = (1; 0.1; \exp(1))$.
- Covariance function $L(\beta, r)$: Gaussian function.
- Classifier: random forest.

Result:

	N	Category	Code name	3x3xH	Support code: S	Multi-Data 127	Support code: S	Support code: S	Support code: S	log-Gaussian Cox
UIUC	1	3122314	4538310	6725275	3722253	6114250	7253245	7666165	7666165	76.31±1.75
	10	979135	706172	8745	344172	3206165	9574115	8423678	95.42±0.71	
	20	679235	801233	979235	3182138	3355131	9373335	9724125	97.80±0.32	
Brodatz	1	655272	754273	883455	3422163	3255091	80334135	7773161	88.91±0.92	
	3	854241	882246	9273091	3173037	3241679	9434045	9273035	95.14±0.41	
	5	823241	828242	8422391	7172253	7262245	7974235	8134153	81.72±1.15	
KTH-Tips	20	826235	804567	8745	347448	3716153	9024131	9215126	92.42±1.11	
	40	897235	888617	903065	3115115	3125097	9433337	9505041	95.74±0.45	

Classification rates and standard deviations of proposed descriptor over 50 random selections compared to state-of-the-art approaches.

N	Code name	PhiSD	PhiSD	PhiSD	PhiSD	PhiSD	PhiSD
1	757416	757416	7666161	7666161	7666161	7666161	7666161
5	3199113	9172411	5142123	9234149	9165414	9165414	9165414
10	354207	9515251	9523065	9541088	9534272	9534272	9534272
15	9534168	9514263	9514263	9575055	9575055	9575055	9575055
20	9724125	9724125	9724125	9724125	9724125	9724125	9724125

Comparison performance of proposed model with the different detector-descriptor types on UIUC dataset.

References:

[1] J. Zhang et al. "Local features and kernels for classification of texture and object categories: a comprehensive study". IJCV, 73(2), p.213-238, 2007.
 [2] G. Csurka et al. "Visual categorization with bags of keypoints". ECCV, p.1-22, 2004.
 [3] H-G. Nguyen et al. "Spatial statistics of visual keypoints for texture recognition". ECCV, pp. 764-777, 2010.
 [4] J. Møller et al. "Log-Gaussian cox processes". SJS, 25(3), p. 451-482, 1998.

RANDOM WALK MODELS FOR GEOMETRY-DRIVEN IMAGE SUPER-RESOLUTION

APPLICATION TO REMOTELY SENSED GEOPHYSICAL FIELDS AT OCEAN SURFACE

Ronan Fabelt

Professor, Institut Mines-Télécom/Telecom Bretagne
Brest, France



Co-authors

Brathim Bousidli
Institut Mines-Télécom/Telecom Bretagne
LabSTIC/TOMS
Brest, France

Emanuelle Autret, Bertrand Chapron
Ifremer/LOS (Space oceanography)
CERSAT
Brest, France

References

- [1] Coaling, E.A., Plank, M.J., Benhamou, S. Random walk models in biology. J. R. Soc. Interface, 5:813-824, 2008.
- [2] Fabelt, R., Pujolle, S., Chessel, A., Benzinou, A., Cao F. 2D image-based reconstruction of shape deformation of biological structures using a levelset representation. CVU, 111(3):295-306, 2008.
- [3] Freeman, W.T., Liu, C. Markov Random Fields for Super-resolution and Texture Synthesis. In A. Blake, P. Kohli, and C. Rother, eds. Advances in Markov Random Fields for Vision and Image Processing, Chapter 10. MIT Press, 2011.
- [4] Galerne, B., Gousseau, Y., Morel, J., 2011. Random phase textures: Theory and synthesis. IEEE Transactions on Image Processing, 20(1), 257-267.
- [5] Monasse, P., Guichard, F. Fast Computation of a Contrast Invariant Image Representation. IEEE TIP, 9:860-872, 2000.



ABSTRACT

This paper addresses stochastic geometry-driven image models and their application to super-resolution issues for textured geophysical fields. Whereas most stochastic image models rely on some priors on the distribution of grey-level configurations (e.g., patch-based models, Markov priors, multiplicative cascades...), we here focus on geometric priors. Regarding image level-lines as realizations of 2D random walks, we introduce a stochastic geometry-driven model for 2D textures and consider an application to image super-resolution. The targeted application is the stochastic interpolation of missing data in multi-sensor sea surface observation.

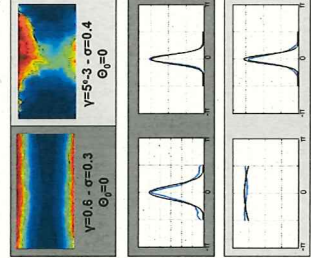
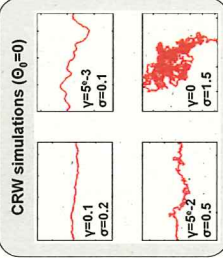
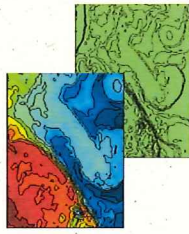


Image level-lines as realizations of 2D random walks

Why level-lines?

Image level-lines provide a contrast-invariant image representation [5] and fully characterize the geometry of an image.



Correlated random walk model [1]:

$$d\theta(s) = -\gamma(\theta(s) - \theta_0) + \sigma dW(s) \quad (1)$$

Directional drift γ Brownian process

Parameters γ and σ control the regularity of the random walk, in terms of regularity along the walk and of oscillation around the directional drift. The CRW model is associated with a Fokker-Planck representation, from which one can derive the stationary statistics of the random direction θ and turning angle $\delta\theta$ [1]:

$$p(\theta) \propto \exp\left(-\frac{\gamma}{\sigma^2}(\theta - \theta_0)^2\right)$$

$$p(\delta\theta) \propto \exp\left(-\alpha \frac{\gamma}{\sigma} \delta\theta^2\right) \quad \text{with } \alpha = \Delta s (1 - \exp(-\gamma \Delta s))$$

Stochastic geometry-driven texture model

Stochastic geometry field model as a generalization of (1) for orientation fields:

$$\nabla \theta = -\gamma(\theta - \theta_0) + \sigma \nabla W \quad (2)$$

Reference orientation θ_0 2D Brownian sheet

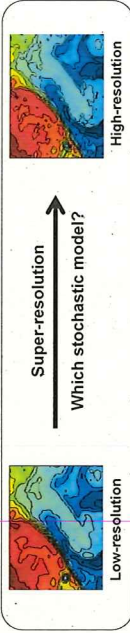
Where $\nabla \theta$ is the gradient in direction θ

Image I such that its level-lines are everywhere tangent to random vector field u_θ :

$$(\nabla I(p), u_\theta(p)) = 0, \forall p \quad (3)$$

Implementation:

- 1) Simulate an orientation field from a numerical integration of the stochastic equation,
- 2) Solve for (3) as a variational minimization (cf. [3]).



Application to texture-based super-resolution of geophysical fields

Problem statement

I_{LR} : High-resolution image
NxM grid

I_{LR} : Low-resolution image
N'xM'K grid (here, K=2^4)

Stochastic super-resolution model

Given a low-resolution image I_{LR} , sample a high-resolution image such that:

$$\left\{ \begin{aligned} \nabla \theta(p) &= -\gamma(p)(\theta(p) - \theta_{LR}(p)) dp + \sigma(p) \nabla W(p) \\ \tilde{I} &= \arg \min_I \int \|\nabla I(p) \cdot u_\theta(p)\| dp \quad \text{subject to } I_{LR} = \mathcal{P}[\tilde{I}] \end{aligned} \right.$$

with W a Brownian sheet and θ_{LR} the orientation field of the low-resolution image (i.e. the angle of the local tangent to the level-lines). \mathcal{P} is an orthogonal projection operator (here a wavelet operator) such that:

$$I_{LR} = \mathcal{P}(I_{HR}) \quad \text{and} \quad \mathcal{P}[I_{HR} - \mathcal{P}(I_{HR})]$$

The low-resolution gradient drives the high-resolution geometrical variability

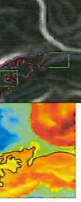
Parameter fields γ and σ are set according to the following observation:

• Large gradients in I_{LR} result in more regular level-lines in I_{HR} .

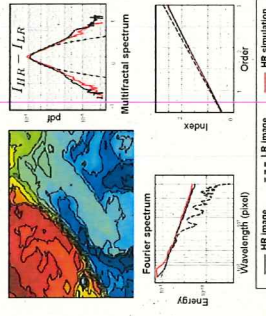
• Conversely, weak gradients in I_{LR} involve regular level-lines in I_{HR} .

$$\gamma(p) = \gamma_0 \|\nabla I_{LR}(p)\|^v \quad \sigma(p) = \sigma_0 \|\nabla I_{LR}(p)\|^{-\beta}$$

(Here, we set empirically $v=2.14, \beta=0.54, \gamma_0=0.14, \sigma_0=0.13$)



Simulation for sea surface temperature observations



Comparison to a Gaussian field simulation: HR image (A), LR image (B), super-resolution with the proposed model (C), and a Gaussian field with the same Fourier spectrum [5] (D).

Contact

ronan.fabelt@telecom-bretagne.eu

Webpage: perso.telecom-bretagne.eu/romanfabelt

Parties prenantes



Auteurs

Yasser Fadlallah

Abdeldjalil Aïssa-El-Bey

Karine Amis

Dominique Pastor

Ramesh Pyndiah

Partenaires



Introduction

- Maximum Likelihood joint detection enables to detect at once the symbols transmitted in the same time intervals.
- ML detector selects the closest point to the received signal in the receive constellation.
- ML is optimal for medium to high SNR values.
- Computation cost increases exponentially with the signal dimension.
- Alternative solution such as sphere decoder keeps an exponential increase of the computation cost.
- **Goal:** find out an efficient detection problem of the transmitted symbols with moderate computation cost.

System Model

- MIMO flat fading channel.
- Perfect knowledge of the channel matrix at the receiver.
- Transmit symbols belong to a finite alphabet constellation.
- Received signal is defined as

Sparse Decomposition

- Let $Q = \{q_1, \dots, q_M\}$ the finite alphabet transmit constellation. Let $\mathbf{q} = [q_1, \dots, q_M]$ the vector space in which the finite alphabet vector can be cast, and \mathbf{B}_q the decomposition matrix defined as

$$\mathbf{B}_q = \begin{pmatrix} \mathbf{q} & \mathbf{0}_M & \dots & \mathbf{0}_M \\ \mathbf{0}_M & \mathbf{q} & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0}_M \\ \mathbf{0}_M & \dots & \mathbf{0}_M & \mathbf{q} \end{pmatrix}$$

- The transmitted symbol vector \mathbf{x} can be rewritten after symbol decomposition on \mathbf{q} as : $\mathbf{x} = \mathbf{B}_q \mathbf{s}$

- The received signal is reformulated as: $\mathbf{y} = \mathbf{H}\mathbf{B}_q \mathbf{s} + \mathbf{z}$.

New MIMO Detector based on ℓ_1 -norm minimization

- The new problem is the decoding of the binary source vector \mathbf{S} . To this end, we propose to solve it using the following minimization problem

$$\arg \min_{\mathbf{s} \in \mathbb{R}^{NM \times 1}} \|\mathbf{s}\|_0, \quad \text{subject to } \mathbf{s} \in \left\{ \|\mathbf{y} - \mathbf{H}\mathbf{B}_q \mathbf{s}\| < \varepsilon, \text{ and } \mathbf{B}_1 \mathbf{s} = \mathbf{1}_N \right\}$$

where ε is a constant defined later, and $\mathbf{B}_1 = \begin{pmatrix} 1_M & \mathbf{0}_M & \dots & \mathbf{0}_M \\ \mathbf{0}_M & 1_M & \ddots & \vdots \\ \vdots & \ddots & \ddots & \mathbf{0}_M \\ \mathbf{0}_M & \dots & \mathbf{0}_M & 1_M \end{pmatrix}$

- In the literature of sparse reconstruction, the ℓ_0 -norm can be relaxed by the ℓ_1 -norm minimization, and the problem is resolved using

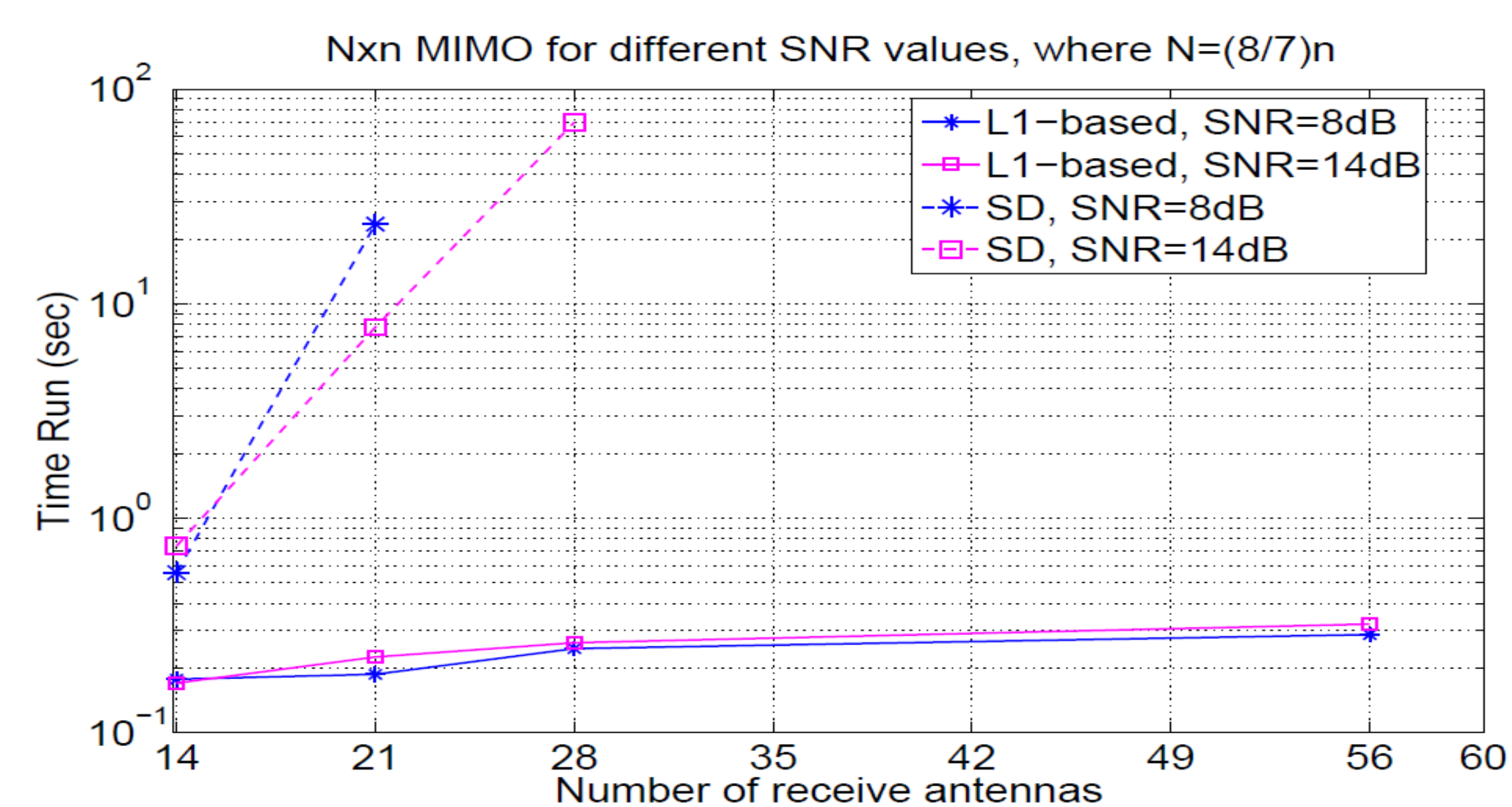
$$\arg \min_{\mathbf{s} \in \mathbb{R}^{NM \times 1}} \|\mathbf{s}\|_1, \quad \text{subject to } \mathbf{s} \in \left\{ \|\mathbf{y} - \mathbf{H}\mathbf{B}_q \mathbf{s}\| < \varepsilon, \text{ and } \mathbf{B}_1 \mathbf{s} = \mathbf{1}_N \right\}$$

Applications

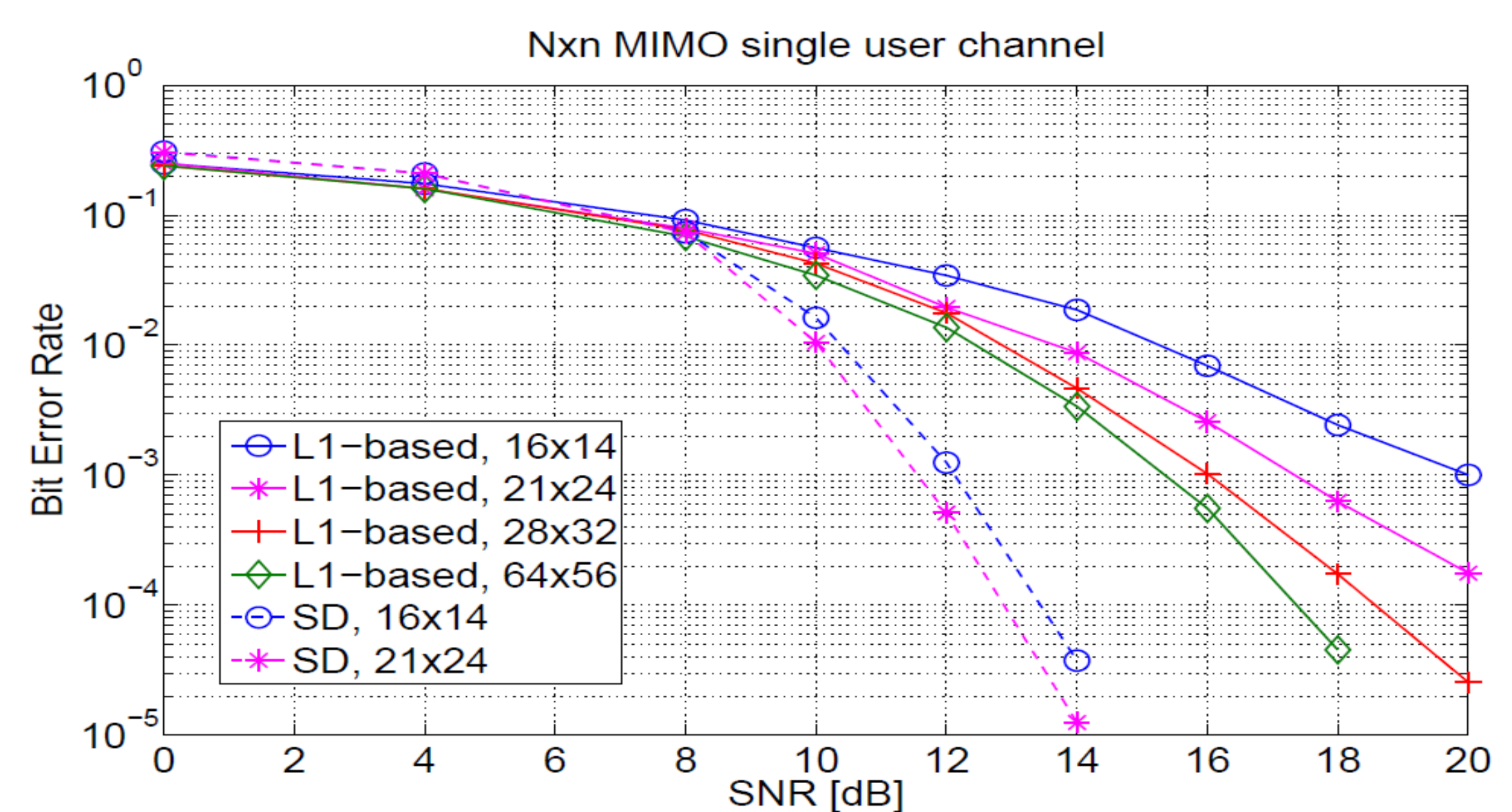
- Large MIMO systems
 - In a noiseless channel, the equivalence between the ℓ_0 -norm and the ℓ_1 -norm hold for large dimensions of \mathbf{s} .
- MIMO frequency selective channel
 - The received signal can be written as

$$\begin{pmatrix} \mathbf{y}(1) \\ \vdots \\ \mathbf{y}(T_f + L) \end{pmatrix} = \begin{pmatrix} \mathbf{H}_0^T & \dots & \mathbf{H}_L^T & \mathbf{0}^T & \dots & \mathbf{0}^T \\ \mathbf{0}^T & \ddots & \ddots & \ddots & \ddots & \vdots \\ \vdots & \ddots & \ddots & \ddots & \ddots & \mathbf{0}^T \\ \mathbf{0}^T & \dots & \mathbf{0}^T & \mathbf{H}_0^T & \dots & \mathbf{H}_L^T \end{pmatrix} \begin{pmatrix} \mathbf{x}(1) \\ \vdots \\ \mathbf{x}(T_f) \end{pmatrix} + \begin{pmatrix} \mathbf{z}(1) \\ \vdots \\ \mathbf{z}(T_f + L) \end{pmatrix}$$

Simulations Results



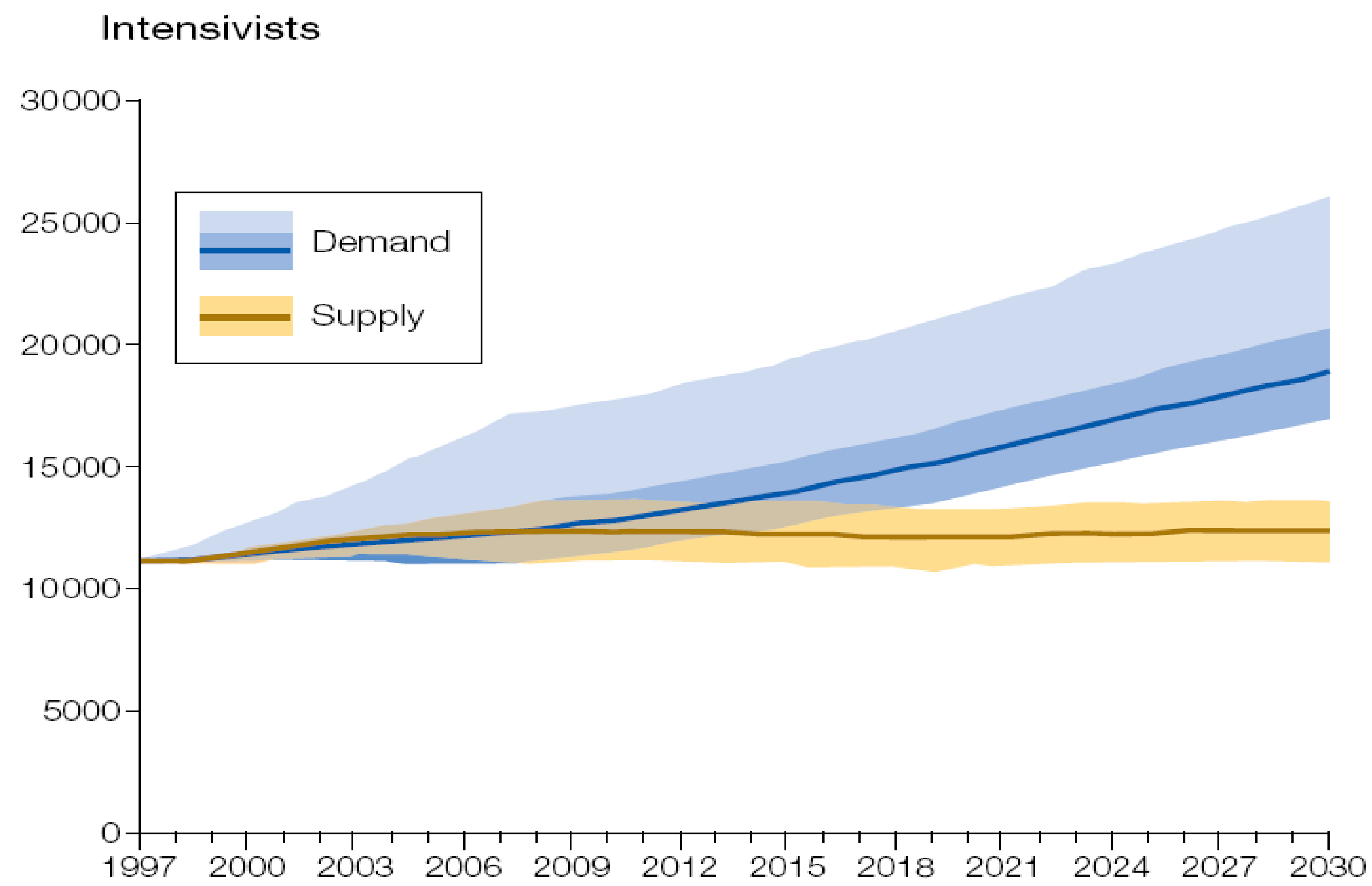
- We assume 4-QAM modulation.
- The computational complexity keeps almost invariant with the system dimensions and the SNR level, whereas the SD time-run increases exponentially with these two factors



Conclusions

- Proposition of a new detection method for determined and underdetermined MIMO systems, based on sparse decomposition of the signal belonging to a finite constellation.
- The proposed detection method is solvable in polynomial time, and uses iterative algorithm such as primal dual interior point method.

Hospital System staffing under pressure



Source: Angus JAMA 2000

The CURVEX solution

Auteurs

Q.-T. Nguyen (TB/Lab-STICC)



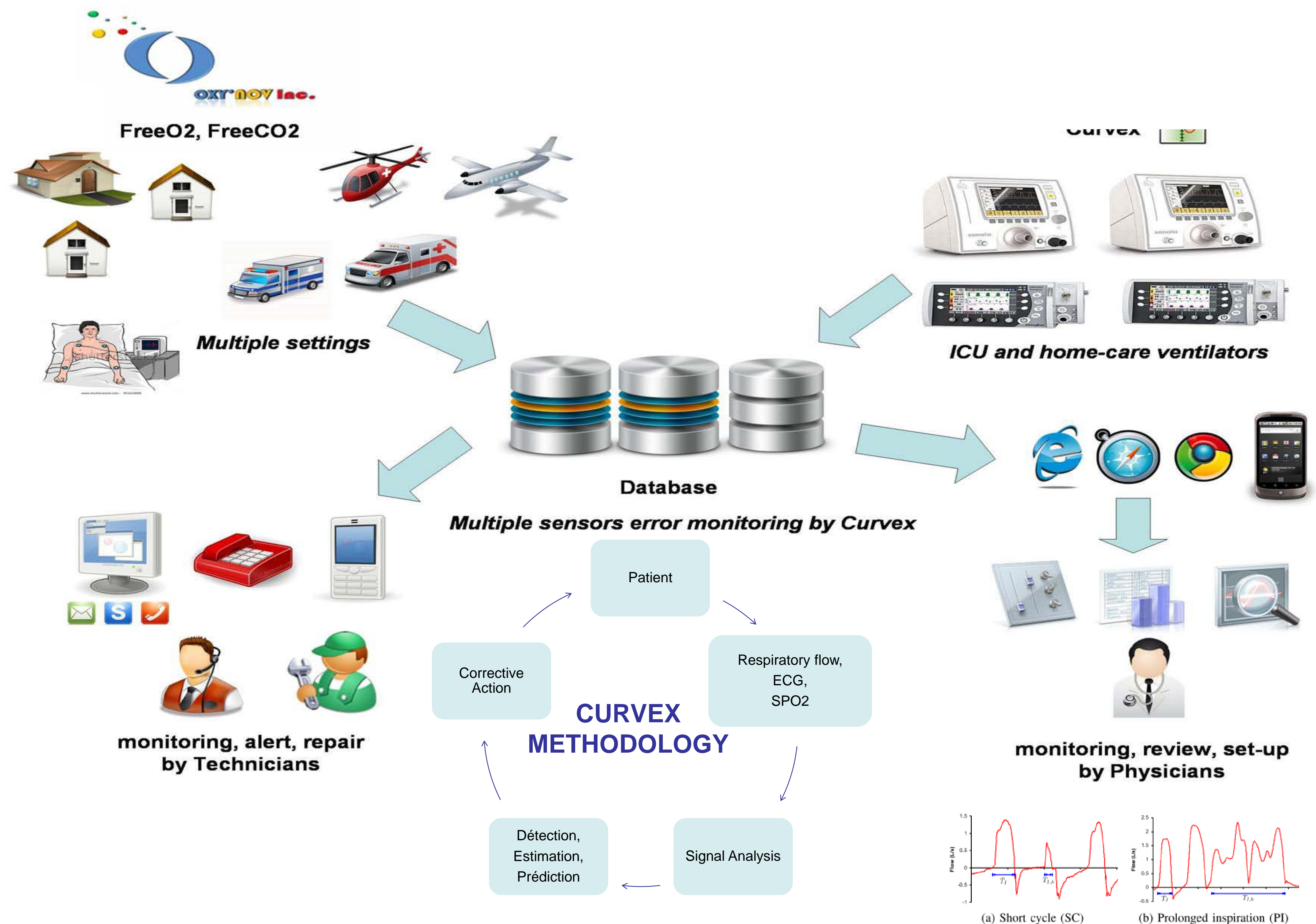
S. Cherif (TB/Lab-STICC)



D. Pastor (TB/Lab-STICC)



E. L'Her (UBO/CHU/LATIM)



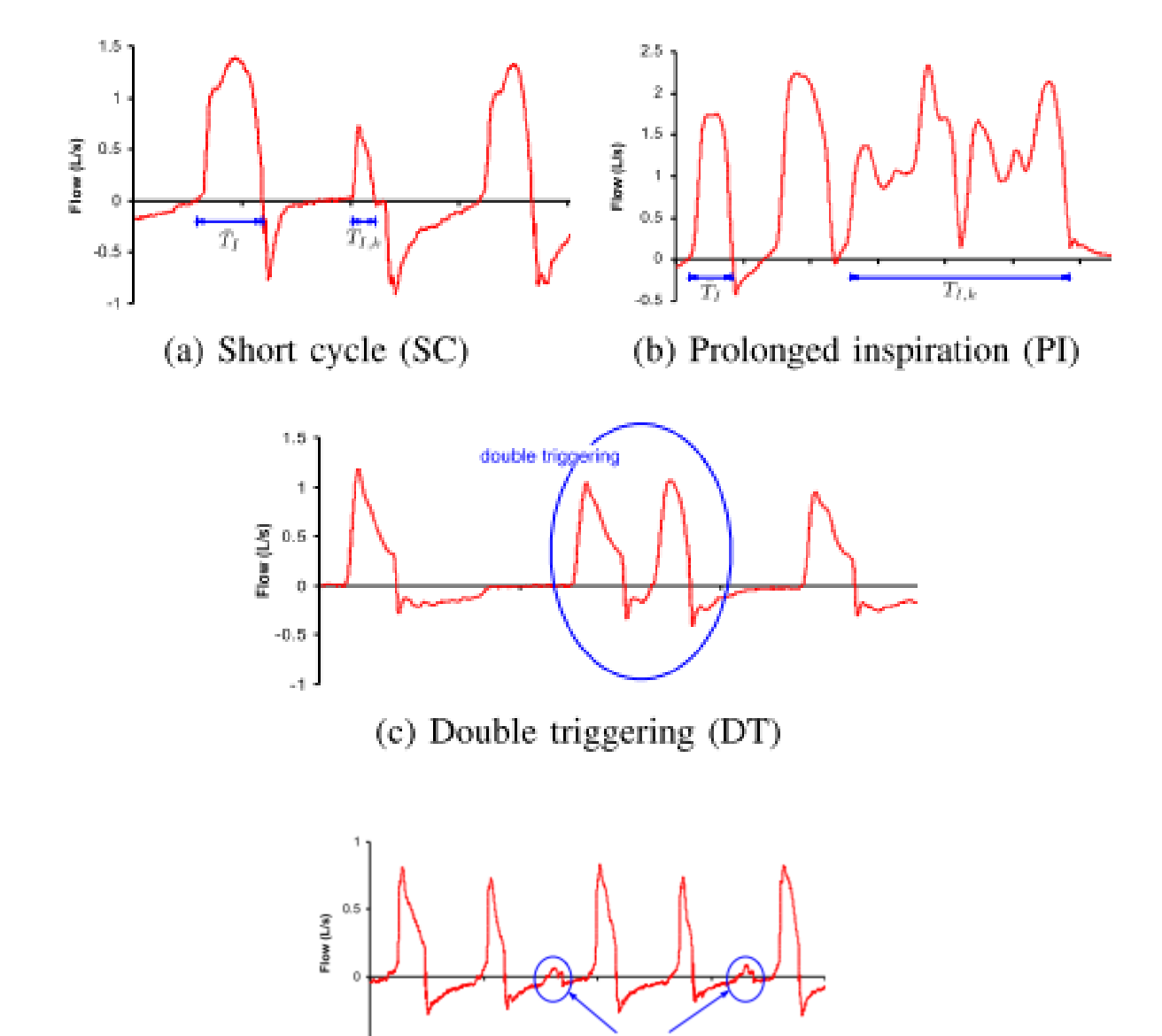
Partenaires



Achievements

- Monitoring of mechanical ventilation (new mathematical framework in robust statistical signal processing, patent FR2988499 - 27/09/2013 « interpretation of expiration curve in mechanical ventilation »)
- Publications in journals, conferences and medicine congresses
- Application to industrial energy management
- Extension to ECG, SPO2, early prediction of patient evolutionary status
- Oxy'nov Inc. (spinoff of Laval University) to open a R&D branch in Brittany

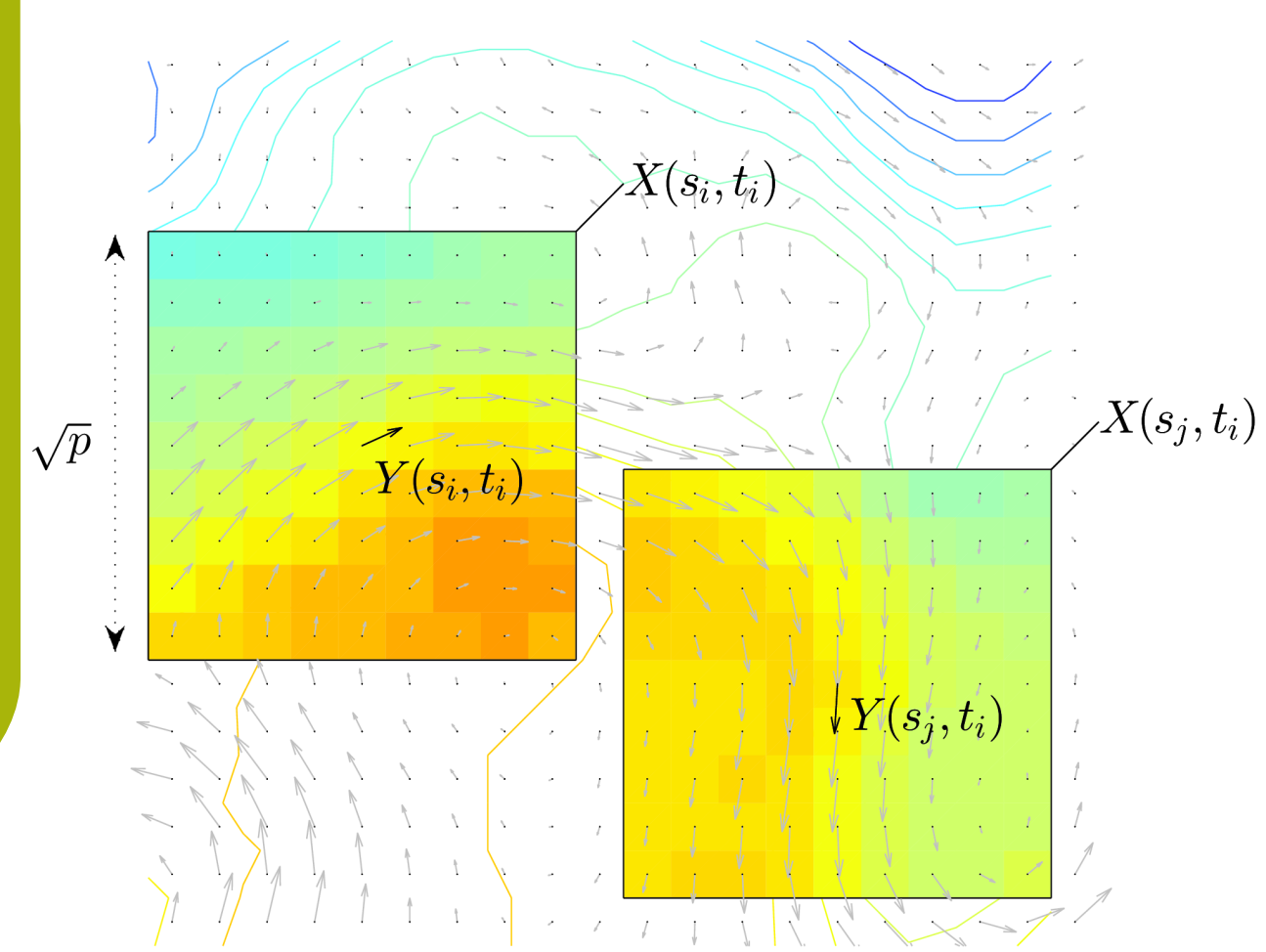
<http://www.telecom-bretagne.eu/lexians/2013/recherche/minute-du-chercheur-dominique-pastor/>



Ex-Vivo analysis of more than 5 000 cycles: works as well in ACV and PSV, either during Invasive or Non-Invasive Ventilation

Accuracy: 93%, Precision 99.5%, Recall 90.5, Specificity 99%

Une approche statistique pour la caractérisation et le suivi des dynamiques superficielles des océans à partir d'images satellitaires



Auteurs

Pierre Tandeo

Sileye Ba

Ronan Fablet

Institut Mines-Telecom

Telecom Bretagne

LabSTICC – TOMS

Brest, France

Bertrand Chapron

Emmanuelle Autret

Ifremer

Laboratoire

d'océanographie spatiale

CERSAT

Brest, France

Partenaires

Ifremer

Lab-STICC

Remerciements



Remote Sensing Systems
www.remss.com

Références

[1] Isern-Fontanet et al. Potential use of microwave sea surface temperatures for the estimation of ocean currents, GRL, vol. 33, pp. L24608, 2006

[2] Lapeyre et al. Dynamics of the upper oceanic layers in terms of surface quasigeostrophy theory, JPO, vol. 36, pp. 165-176, 2006

[3] De Sarbo et al. A maximum likelihood methodology for clusterwise linear regression, Journal of Classification, vol. 5, pp. 249-282, 1988

Résumé

Les mesures satellitaires de courants (U,V) et de température de l'eau de surface (SST), fournissent une information sur les dynamiques de l'océan. Certaines études (cf. [1] et [2]) ont montré que les champs de température peuvent être considérés, dans certaines situations, comme des traceurs actifs de la dynamique de surface. Dans ce cas, de fortes corrélations existent entre les variations locales de SST et les courants (U,V). Existe-t-il d'autres relations entre la température et les courants ? Quand et où la SST peut-elle être considérée comme un traceur actif ou passif ? Dans cette étude, nous mettons en place une méthode statistique et explorons un historique d'observations satellitaires pour identifier et suivre des modes dynamiques cachés.

Méthode

□ K fonctions de transfert cachées entre :

Y → courant (U,V) en un point

X → température SST au voisinage (patch)

$$Y(s_i, t_i) = \sum_{k=1}^K H_k(X(s_i, t_i))$$

□ Identification des K fonctions de transfert à partir d'un modèle de régressions linéaires latentes (cf. [3]) :

$$p(Y|X, \theta) = \sum_{k=1}^K \lambda_k \mathcal{N}_k(Y; X\beta_k, \Sigma_k)$$

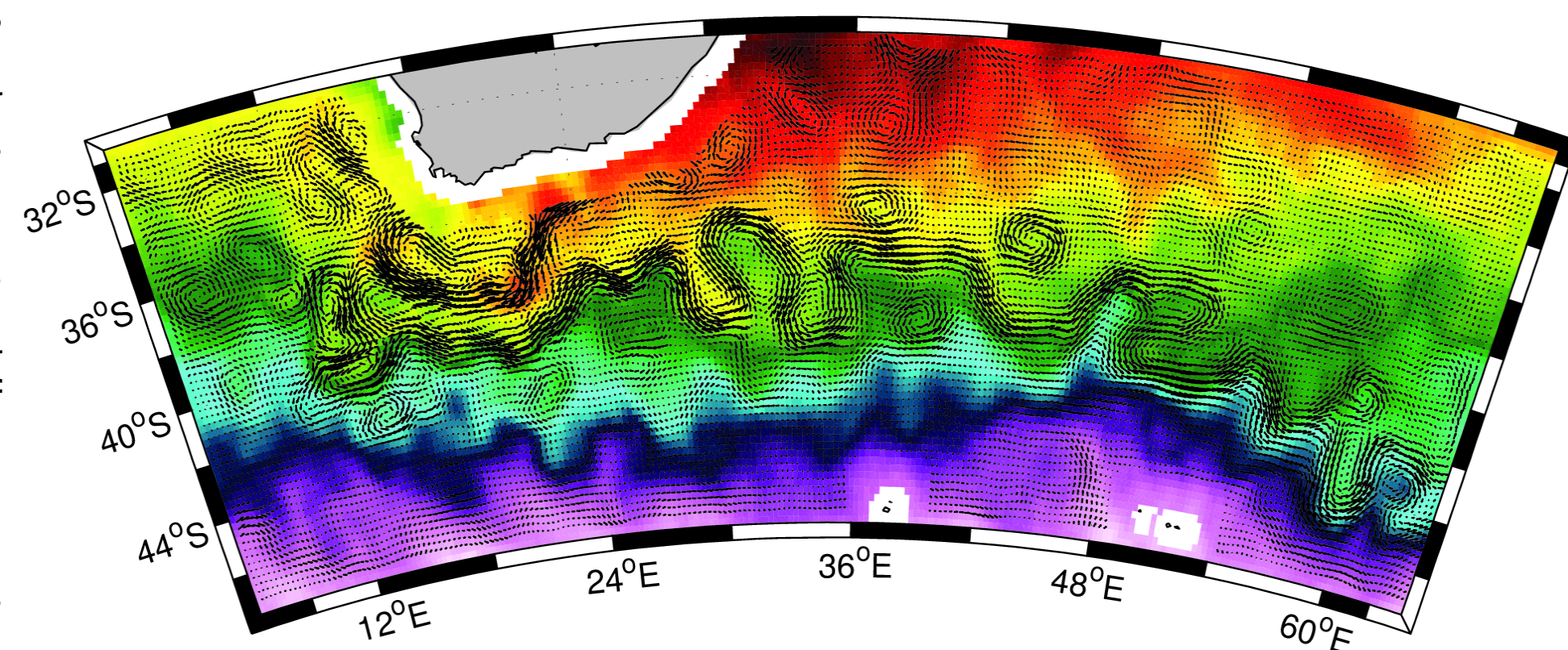
□ Estimation des paramètres par l'algorithme EM

□ Suivi des modes dynamiques à partir des cartes de probabilités *a posteriori*

Données

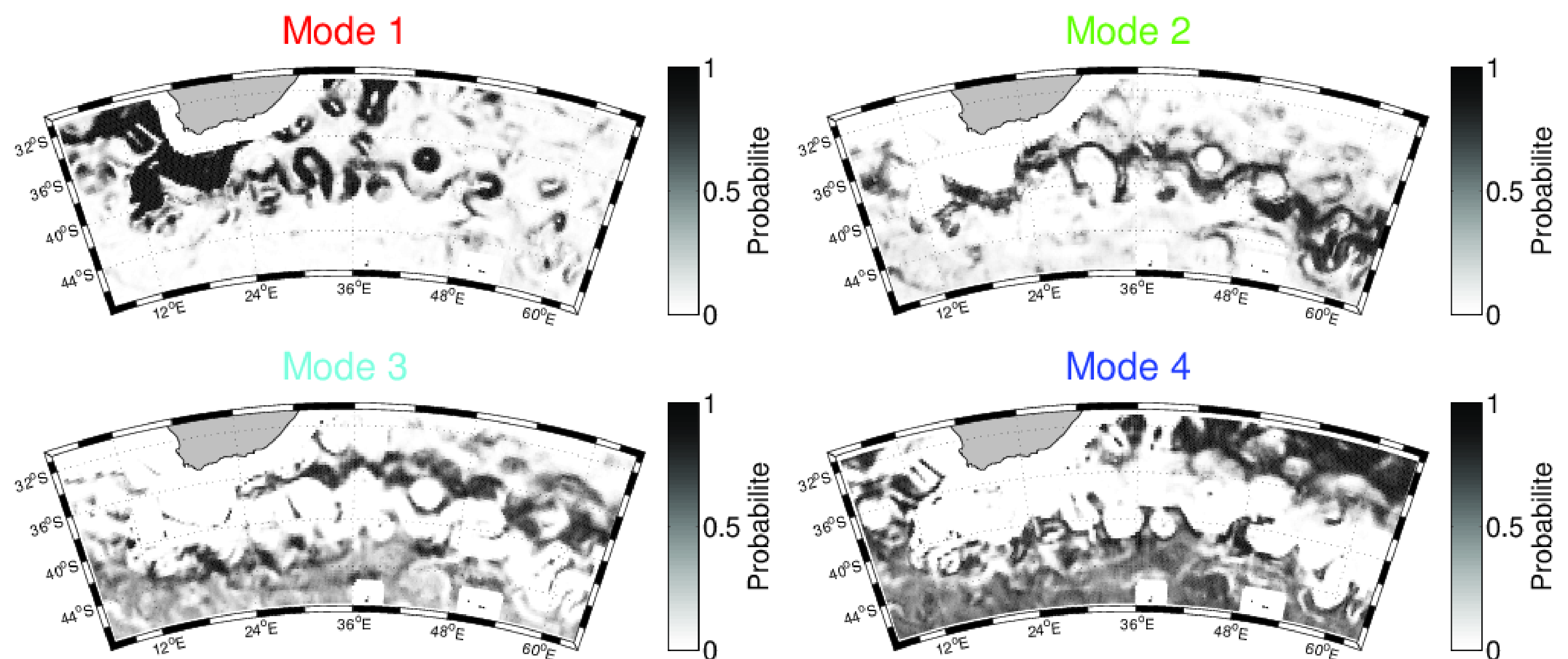
□ Température SST :
Produit interpolé RSS
Journalier au 1/4°

□ Courant (U,V) :
Produit interpolé AVISO
Journalière au 1/3°



Résultats

	Mode 1	Mode 2	Mode 3	Mode 4
Courant (intensité)	fort (principal Aiguilles + tourbillons)	fort (secondaire Aiguilles)	faible	faible
Courant (sens)	Nord-Sud	Ouest-Est	Est-Ouest	Ouest-Est
Température	élevée	élevée	moyenne	froide
SQG-like		✓		✓



Conclusion

La méthode proposée permet un apprentissage à l'aveugle, sans *a priori* physique, de fonctions de transfert cachées entre la SST et (U,V). Ces relations cachées correspondent à différents modes dynamiques dont certains s'apparentent à la théorie SQG. A partir du calcul des probabilités *a posteriori*, nous pouvons suivre l'évolution spatio-temporelle de ces modes dynamiques.

Perspectives

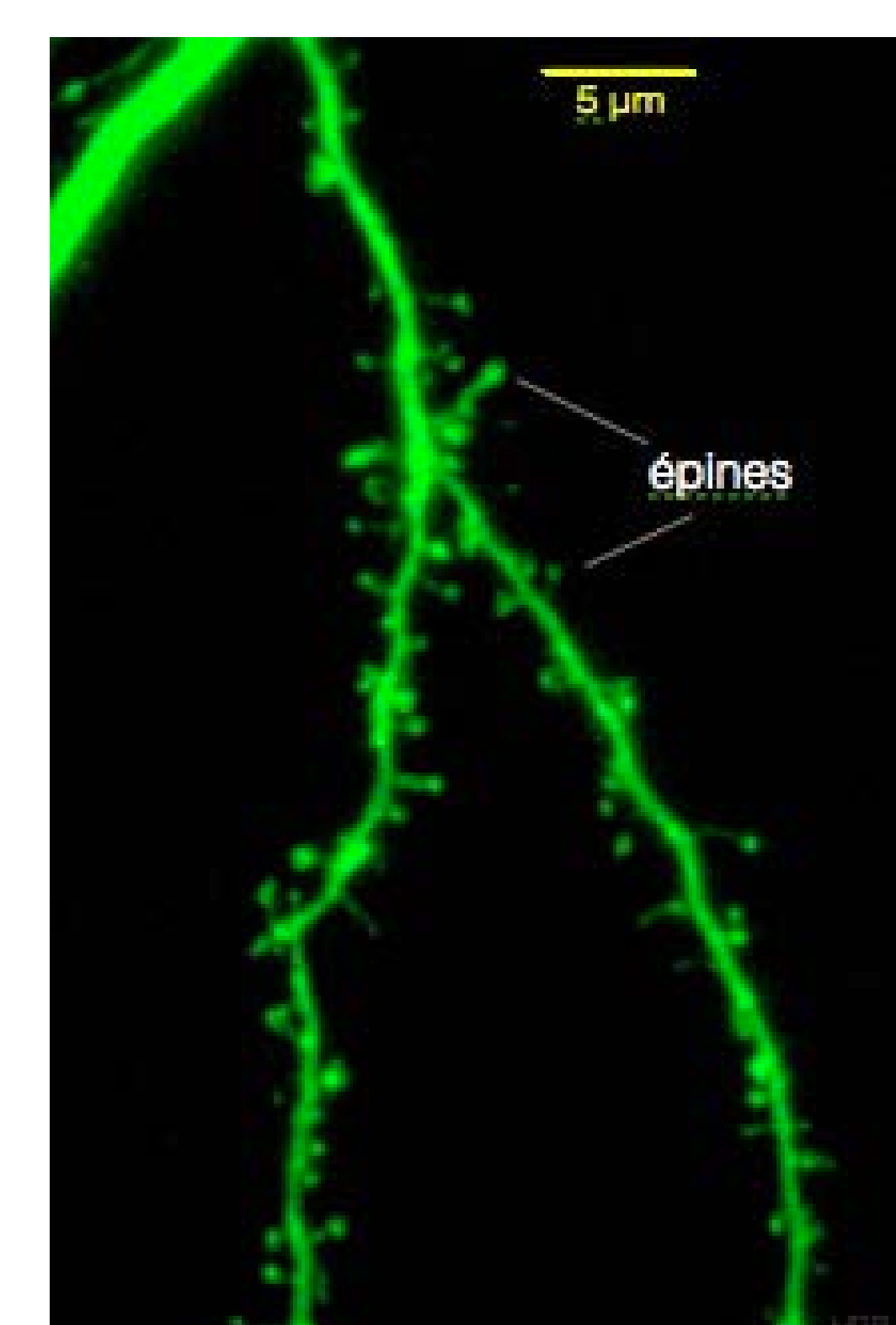
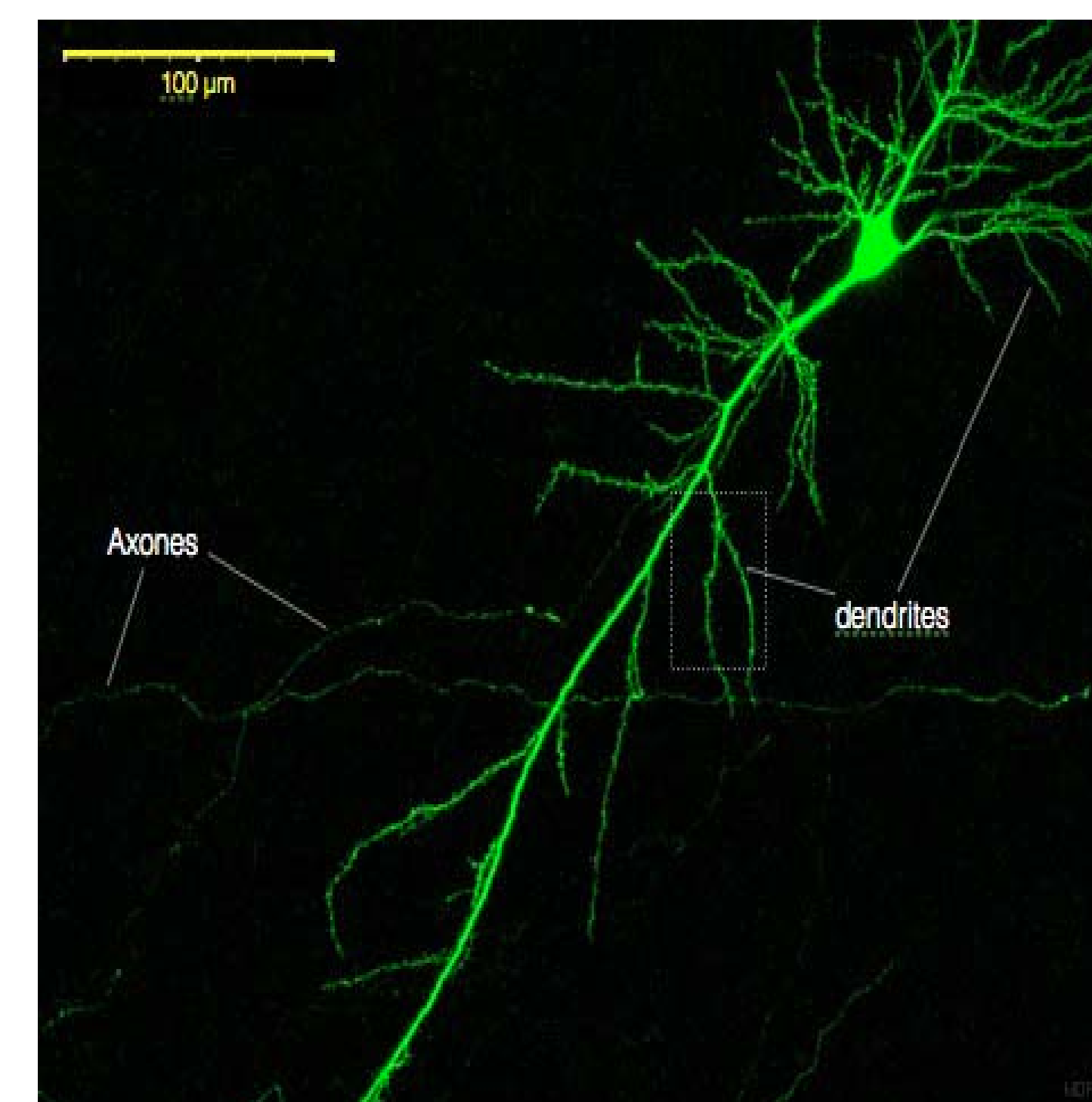
L'utilisation d'autres traceurs actifs tels que la salinité (SSH) ou la couleur de l'eau (Chl-a) ainsi que des données à haute résolution spatiale permettrait un raffinement du modèle. De plus, le suivi des probabilités *a posteriori* par un modèle stochastique et l'utilisation des K fonctions de transfert permet d'envisager une estimation des courants de surface à partir de la température.



DE NOUVEAUX RÔLES POUR UNE STRUCTURE ANCIENNEMENT CONNUE

Les clusters d'épines discrètes agissent comme des unités de calcul

- Les épines dendritiques ont été identifiées anatomiquement par Ramon Y. Cajal en 1911 qui les a qualifiées de "espinas" en raison de leur ressemblance avec des épines sur les tiges des fleurs.
- L'idée que les dendrites ne sont que des câbles passifs qui relaient les signaux entrants sur le corps de la cellule ne tient plus.
- Des études récentes révèlent que les sections dendritiques contenant des clusters d'épines discrètes agissent comme unités de calcul (Blom H, Rönnlund D, Scott L, Westin L, Widengren J, et al., *Nature*, 2013).
- Cette clusterisation est influencée par les entrées sensorielles (Frost N. A., Shroff H., Kong H., Betzig E., Blanpied T. A., *Neuron*, Vol. 67, Issue 1, 15 July 2010).



Plasticité structurale et apprentissage

- Les épines dendritiques subissent en réalité d'une part des changements de leur forme et d'autre part un turn-over permanent (elles apparaissent et disparaissent)
- Lamprecht and Le Doux (2004) proposent une revue des mécanismes de plasticité structurale associées aux épines dendritiques du neurone post-synaptique. Ces mécanismes sont généralement associées au renforcement et à la stabilisation de l'apprentissage synaptique.

Auteurs

Ehsan Sedgh Gooya



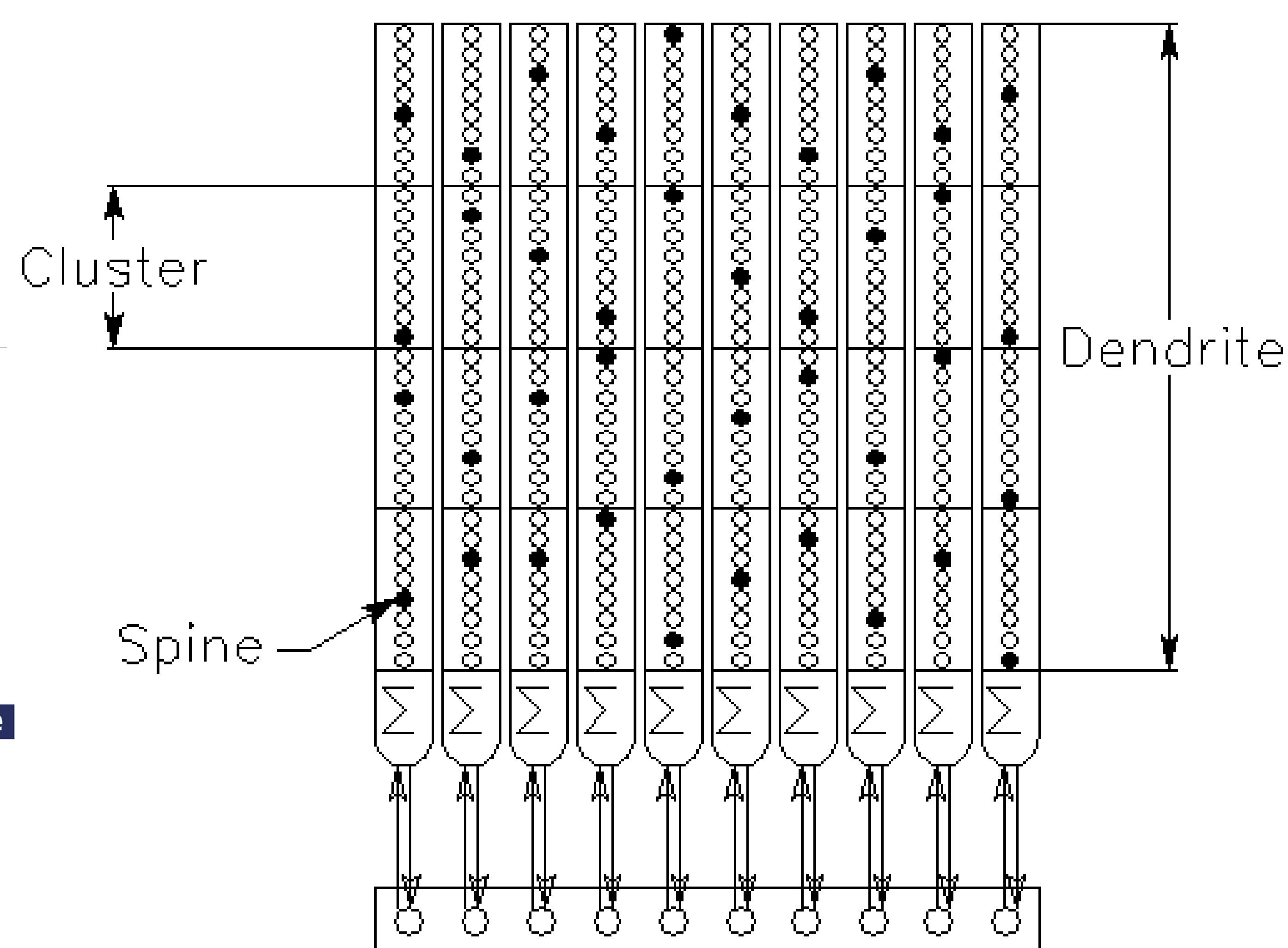
Dominique Pastor



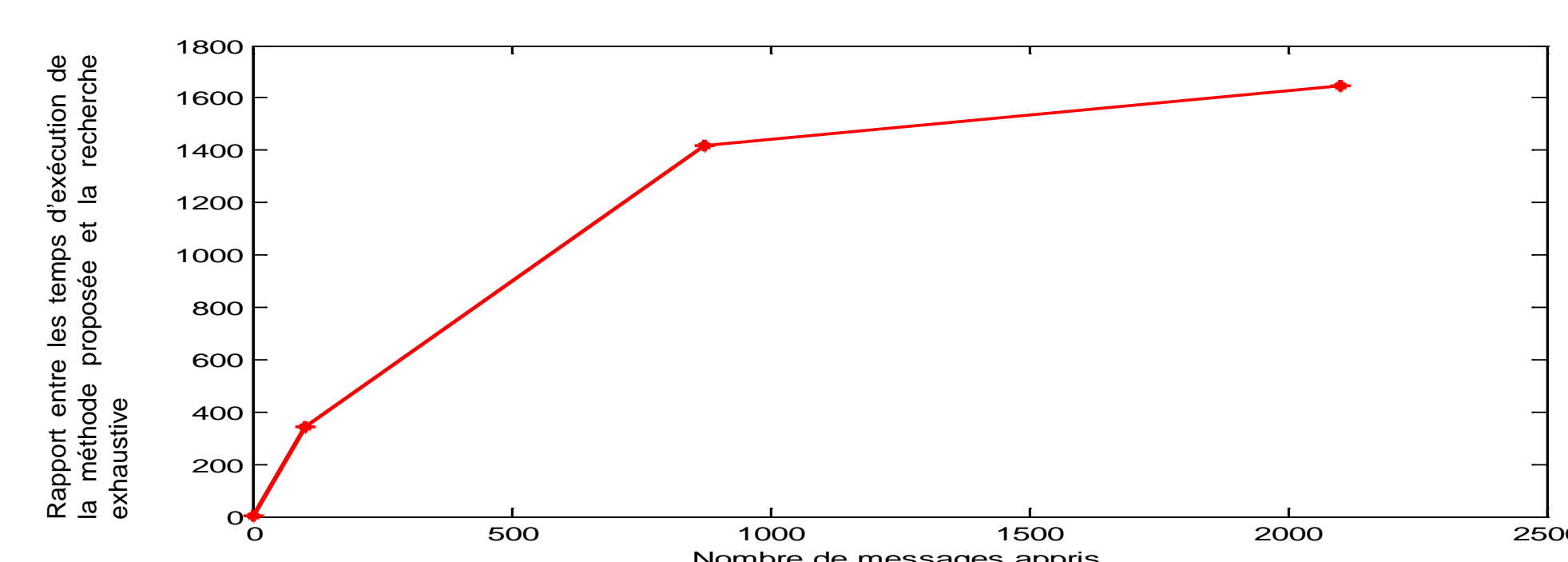
Partenaires



MODELE FORMEL ET PROPRIETES



- Comme les réseaux à cliques de Berrou et Gripon, ce modèle n'emploie que des connexions binaires (pas de poids synaptiques)
- La taille de la mémoire (le stockage) n'est pas un réel problème. Par contre, il faut savoir retrouver très vite et sans erreur un message, ce que permet le neurone formel proposé
- La formalisation mathématique montre que ce neurone ne commet aucune erreur: il reconnaît tout ou une partie des messages qu'il a appris et ne reconnaît aucun message ou partie de message qu'il n'aurait pas appris
- Il est aussi performant qu'une recherche exhaustive, mais pour un coût de calcul largement inférieur.



PERSPECTIVES

- Introduction de la plasticité à l'aide d'un paramètre de tolérance afin de traiter le cas de messages distordus et/ou bruités
- Structure distribuée et/ou hiérarchique (réseau de neurones)
- Concept d'apprentissage attentif
- Vers des machines auto-apprenantes

- Les calculs de cette simulation sont faits en utilisant une programmation Matlab et des signaux générés aléatoirement
- Pour cette simulation, le dendrite est composé de 4 clusters et 8 épines par clusters. Chaque dendrite est associée de manière univoque à un message appris.
- Aucune itération de décodage
- On apprend en ajoutant du matériel, sans altérer ce qui a déjà été appris et en gardant les mêmes capacités de discrimination

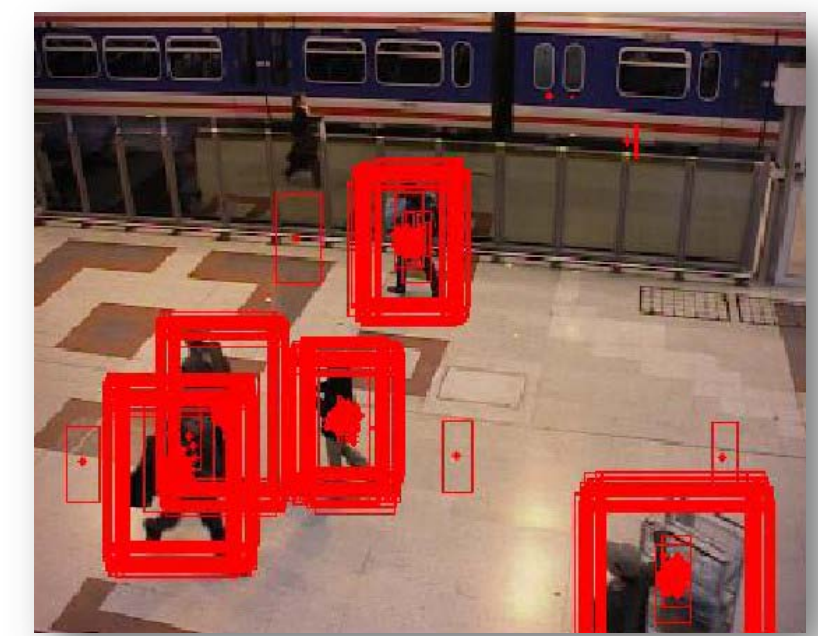
CONTEXTE

- Très grand volume de données vidéo à analyser dans de nombreux domaines d'application : vidéosurveillance, robotique, véhicules autonomes, interactions homme-machine, imagerie du vivant...
- Besoin de systèmes autonomes et intelligents, capables d'extraire automatiquement les informations utiles et de les interpréter.

OBJECTIF, PROBLEMATIQUE ET METHODES

Développer des modèles et des méthodes pour détecter et suivre de façon automatique, fiable et robuste des objets multiples dans des séquences vidéo issues d'une ou de plusieurs caméras.

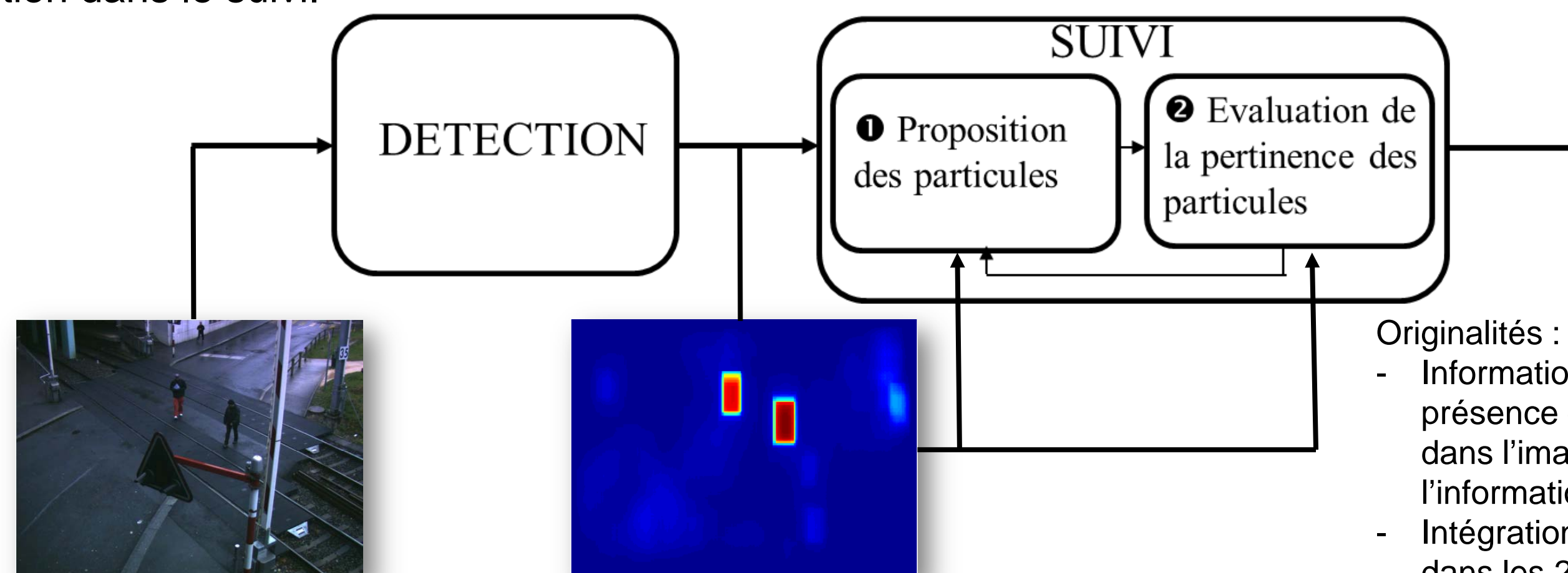
- **Principales difficultés** : nombre croissant d'objets (grande dimension), apparitions / disparitions d'objets, bruit, fausses détections, complexité de l'environnement, non-stationnarités (variations de l'environnement, du mouvement et de l'apparence des objets), occlusions...
- **Problème d'estimation séquentielle** : déterminer le nombre d'objets et leurs paramètres caractéristiques au cours du temps.
- **Outils méthodologiques** :
Méthodes séquentielles de Monte-Carlo (filtrage particulaire),
Méthodes de Monte-Carlo par chaînes de Markov (MCMC).



Exploration de l'espace d'état par les particules

TRAVAUX ACTUELS

- **Optimisation des liens entre détection et suivi** : intégration de la sortie « soft » d'une méthode de détection dans le suivi.



- **Gestion des variations de mouvement et d'apparence des objets** : nouveaux modèles dynamiques.
- **Amélioration du suivi en grande dimension** :
 - Versions séquentielles des méthodes MCMC plus performantes que les filtres particulaires lorsque la dimension augmente,
 - Lois de proposition plus efficaces pour explorer l'espace d'état et guider rapidement les algorithmes vers les zones à forte vraisemblance.



Suivi multi-objets : (1) suivi de 4 personnes, (2) détection automatique d'une 5^{ème} personne, (3) 1^{ère} occlusion partielle, (4) 2^{ème} occlusion partielle, (5) suivi après les occlusions

Parties prenantes



Auteurs

Christelle Garnier
Mehdi Oulad Améziane
François Septier
Yves Delignon
Emmanuel Duflos

Partenaires

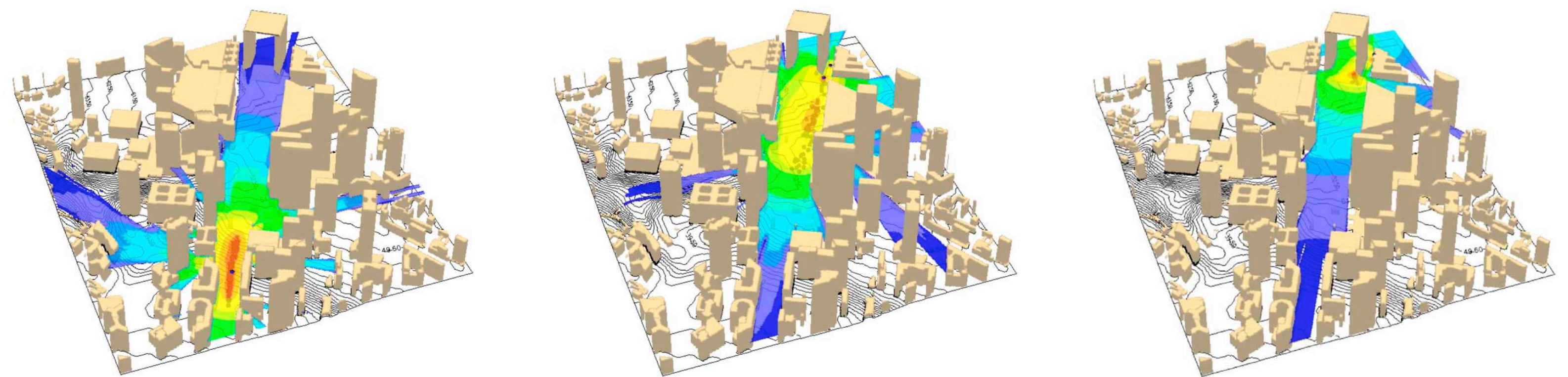


Contexte

- Menace grandissante de rejets délibérés ou accidentels d'agents nucléaires, radioactifs, biologiques ou chimiques (NRBC) ayant des conséquences dramatiques pour la population et l'environnement
- Mise en place d'un réseau mondial de capteurs dans le cadre du Traité d'Interdiction Complète des Essais Nucléaires (TICEN)

Objectifs

Développer une méthode de détection et de localisation de sources de rejets polluants depuis des mesures bruitées de concentration issues de multiples capteurs



Simulation d'un rejet d'agent toxique au coeur du quartier de la Défense à Paris

Problématiques

- Complexité des modèles météorologiques nécessaires pour la simulation réaliste de la dispersion atmosphérique d'agents toxiques.
- Imperfection et inhomogénéité des capteurs utilisés.
- Besoin d'une solution rapide et fiable afin de minimiser les conséquences d'un rejet.

Travaux actuels

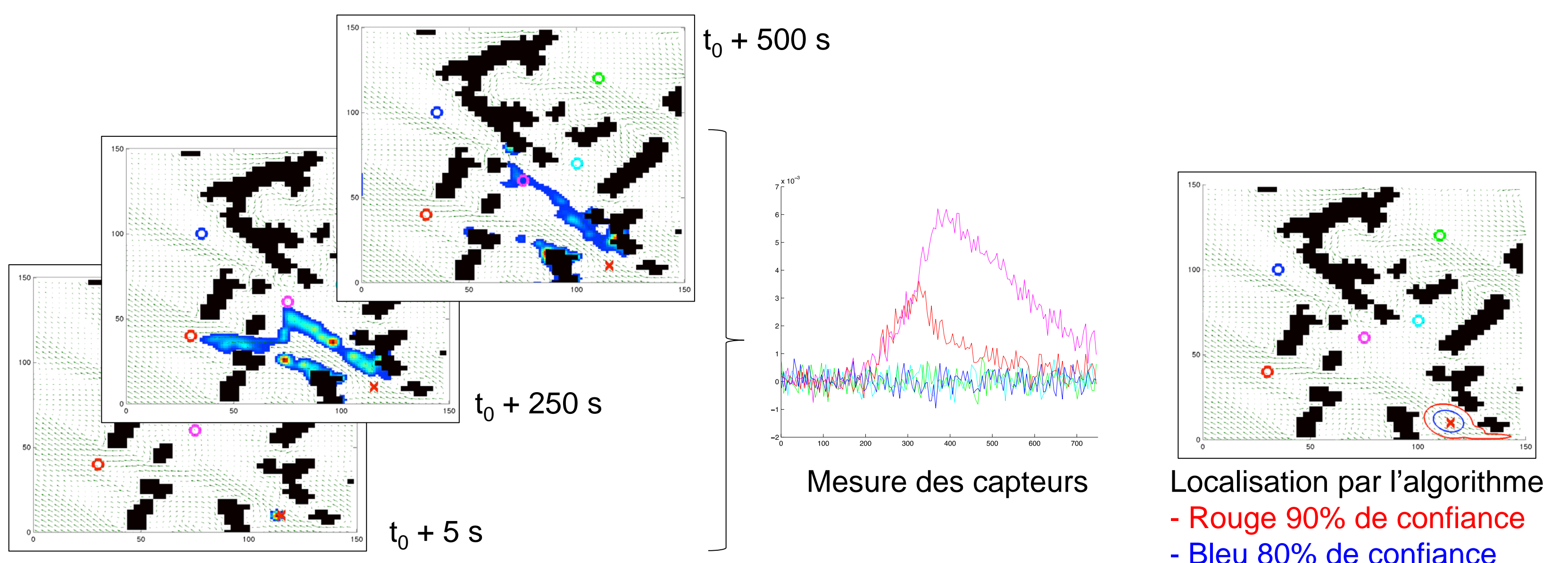
Développement d'une méthode statistique de type Monte-Carlo permettant contrairement aux approches existantes de:

- Fournir un intervalle de confiance sur l'estimation fournie à l'utilisateur.
- Exploiter un modèle de dispersion atmosphérique de flux de turbulences complexes par des modèles Lagrangien (Parallel Micro-SWIFT-SPRAY, PMSS)
- Converger plus rapidement vers les zones de rejet les plus probables grâce à l'utilisation de techniques adaptatives (algorithmes PMC, AMIS).

Partenaires



UMR CNRS 8219



Scénario d'un rejet à t_0 d'une source localisée à (115,10) dans un environnement avec bâtiments (noir) et 5 capteurs (cercles)

Riemannian Geometry for 3D Human Video Retrieval

Context and Issues

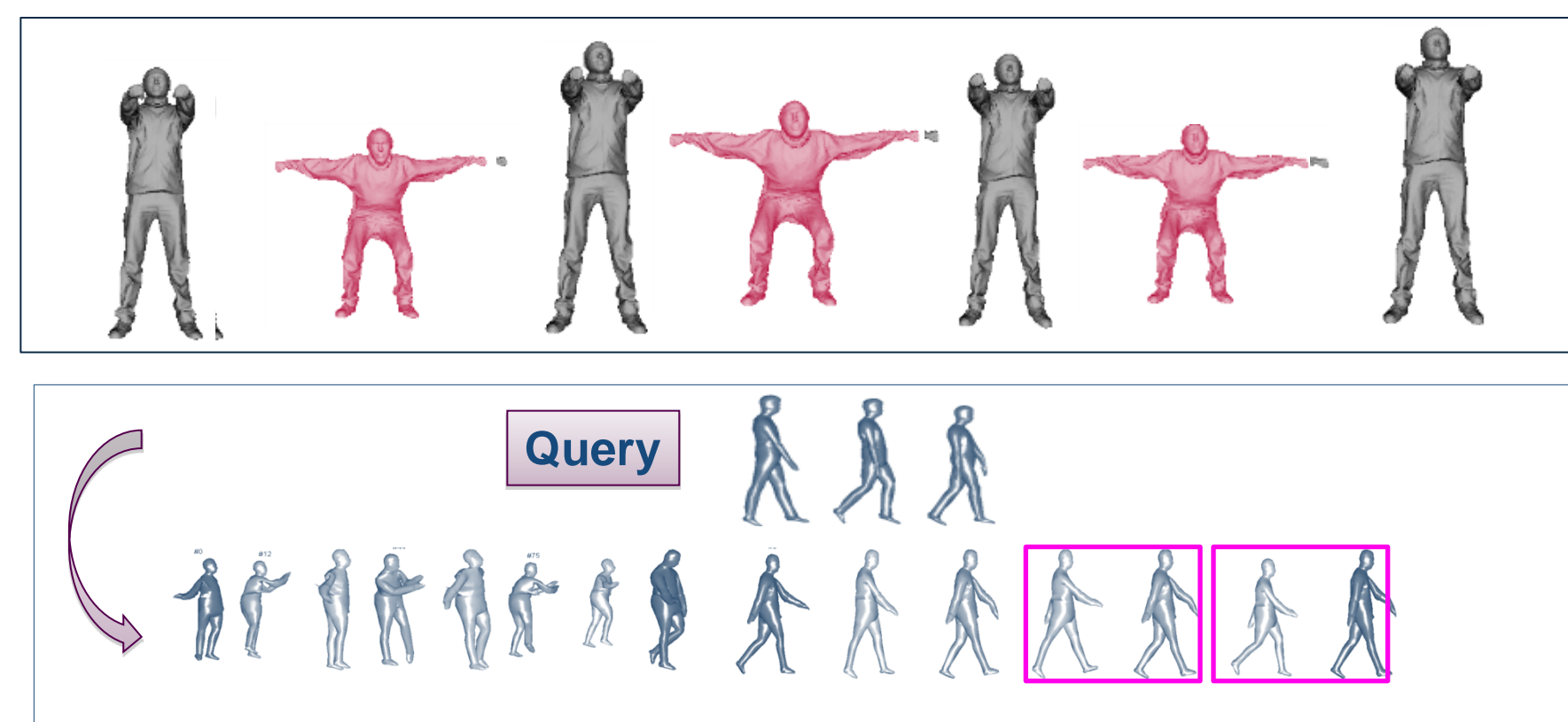
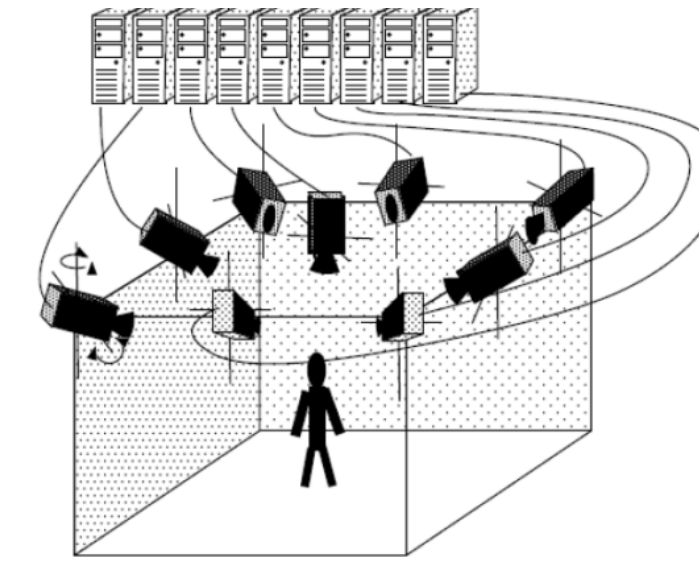
- Long sequence of 3D videos : massive amounts of data

🔍 browsing and searching for relevant information quickly become difficult

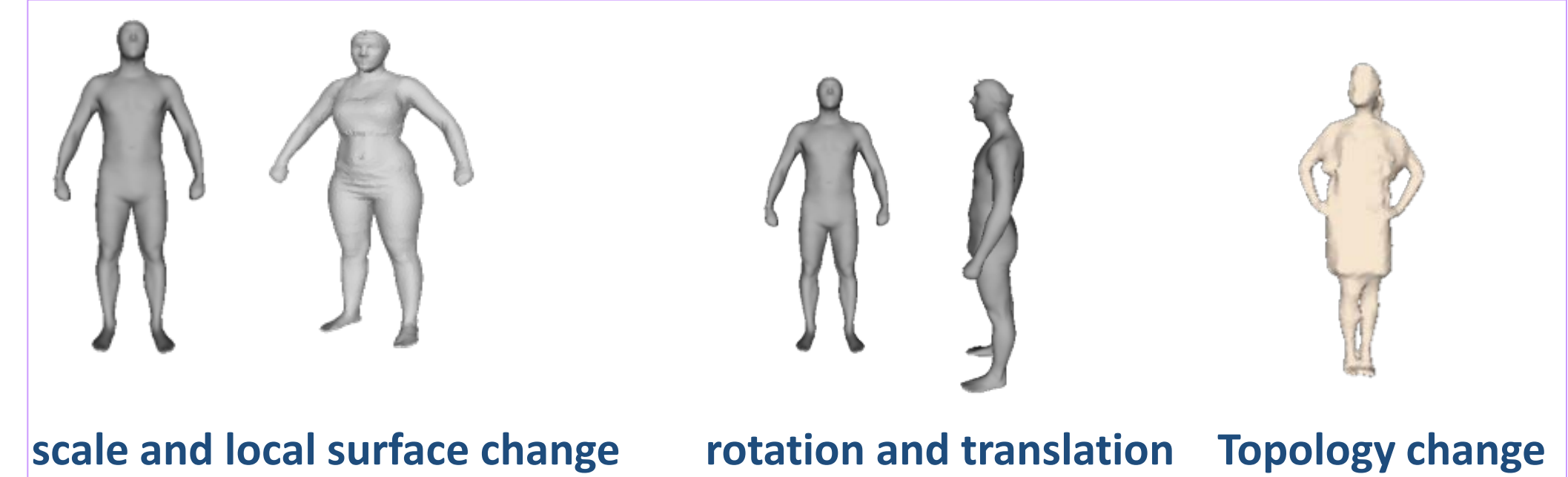
➔ Need for 3D video segmentation system

➔ Pose/Motion retrieval

➔ Video summarization

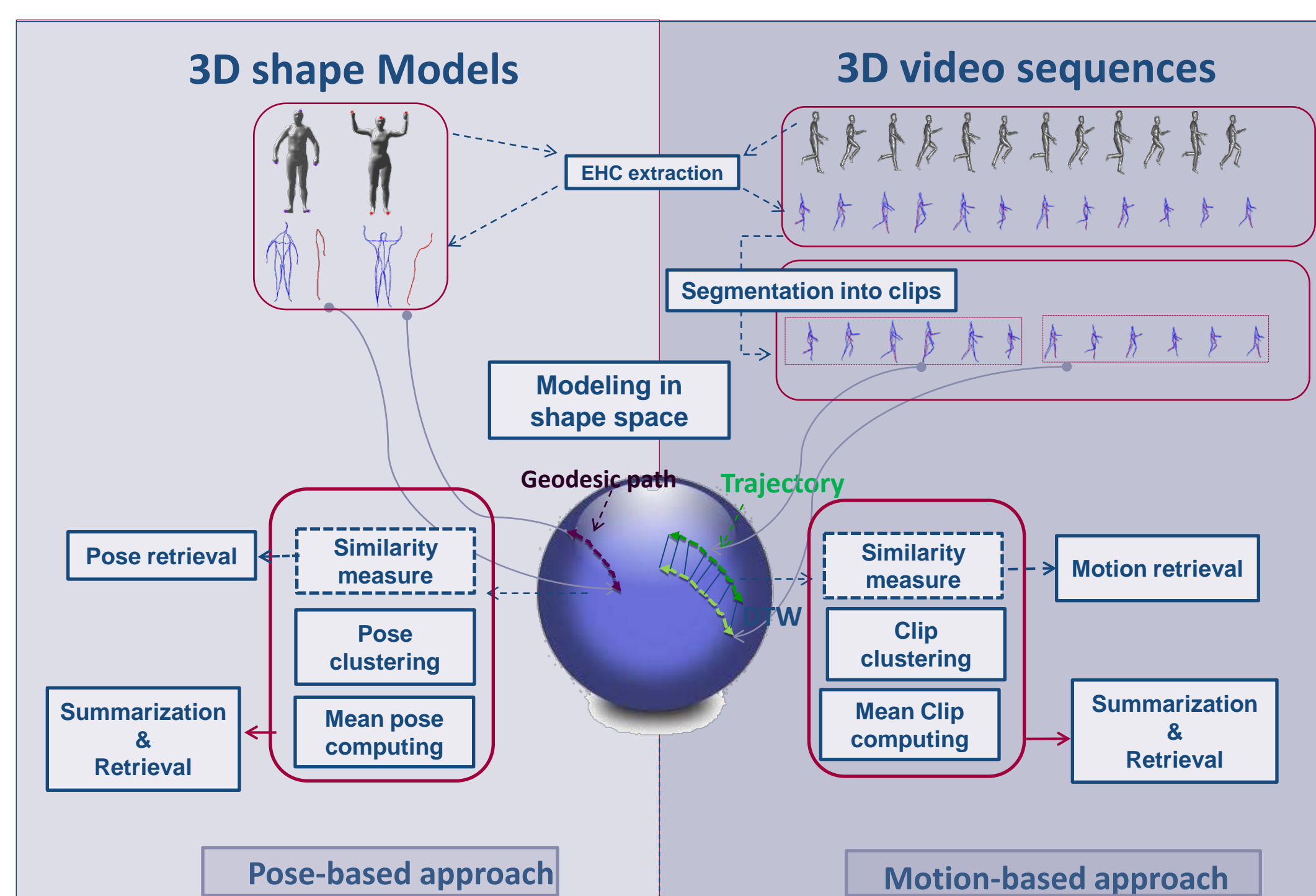


• Challenges:

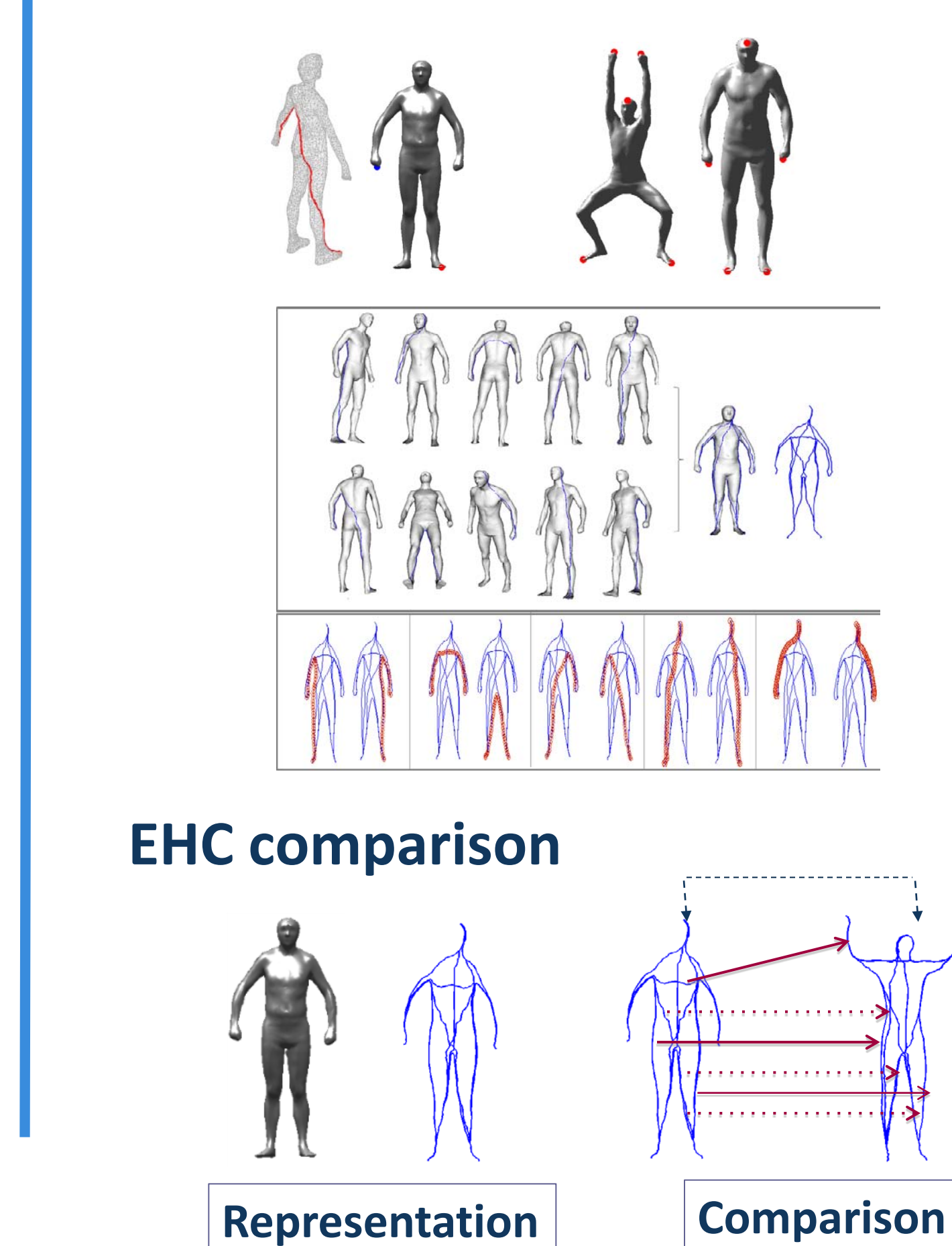


Approach: Riemannian Geometry

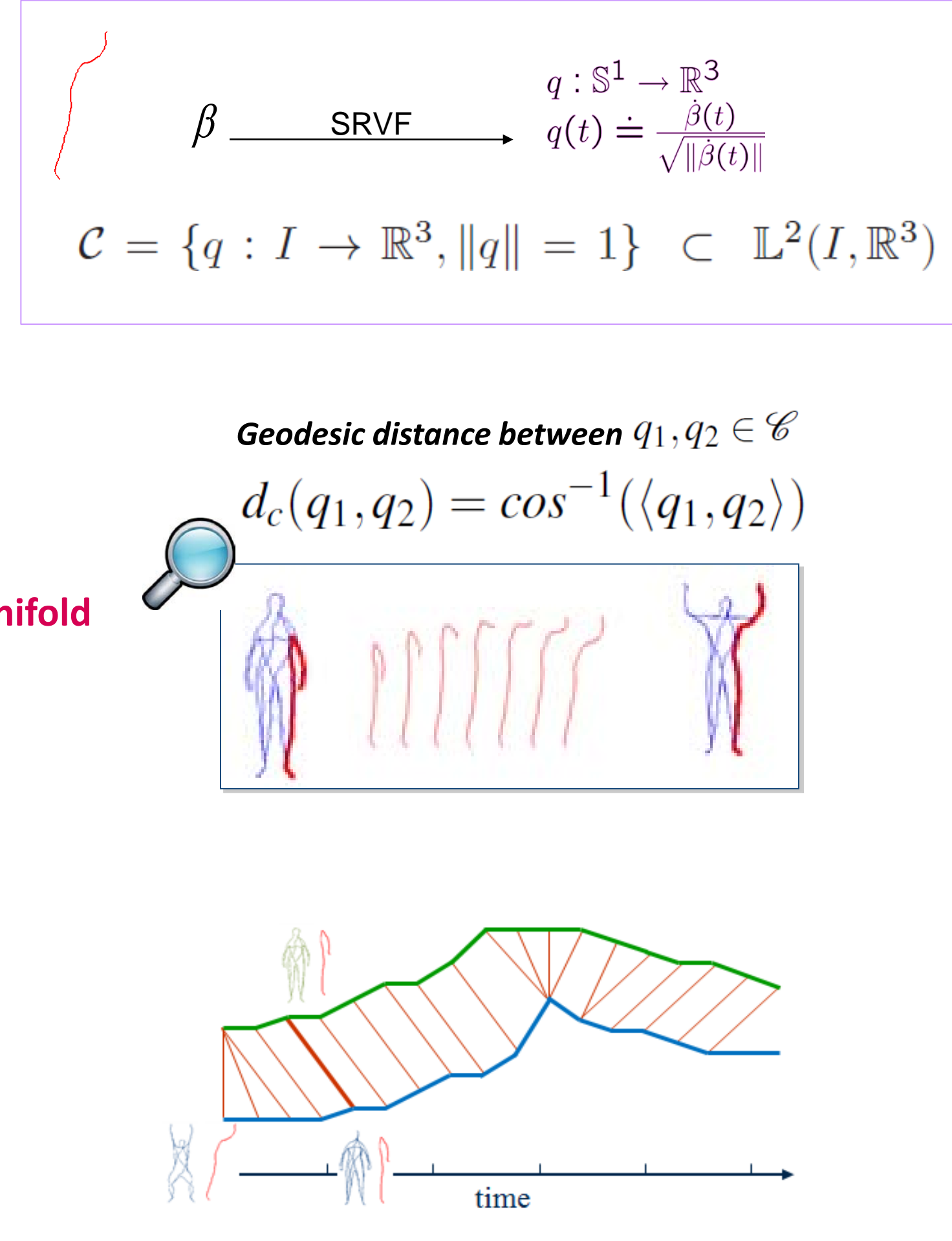
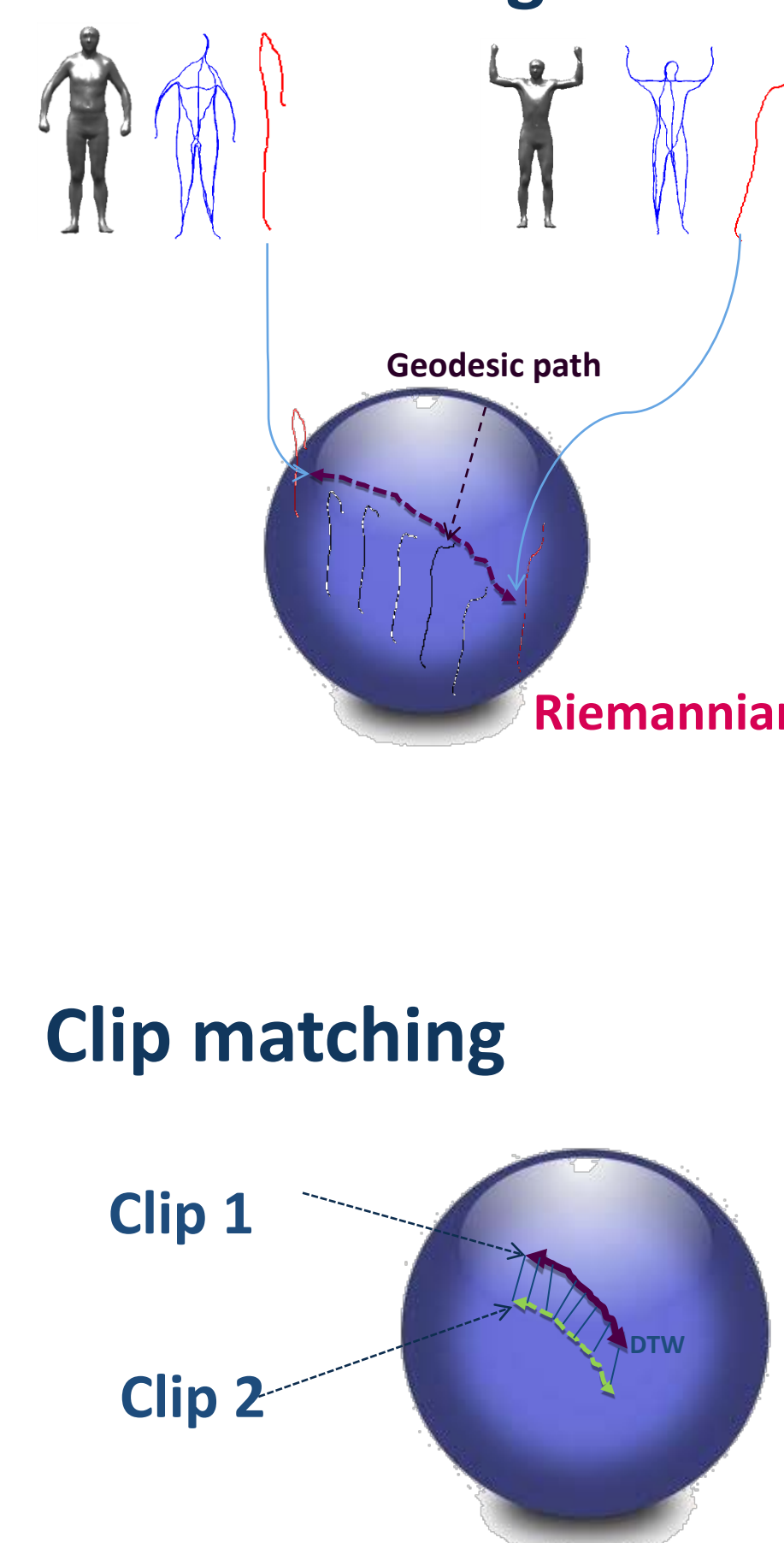
Overview



Extremal Human Curve extraction

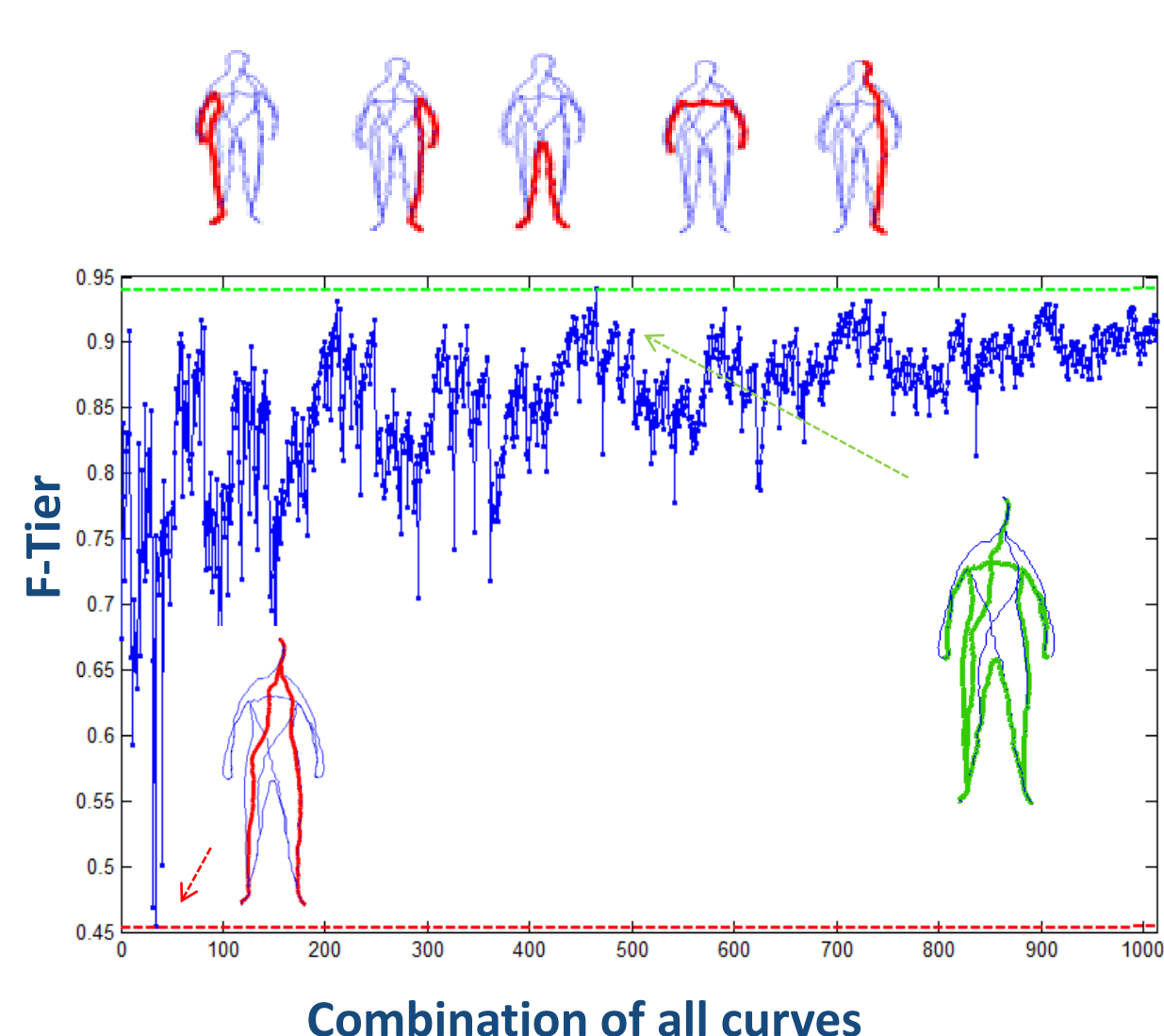


Pose matching

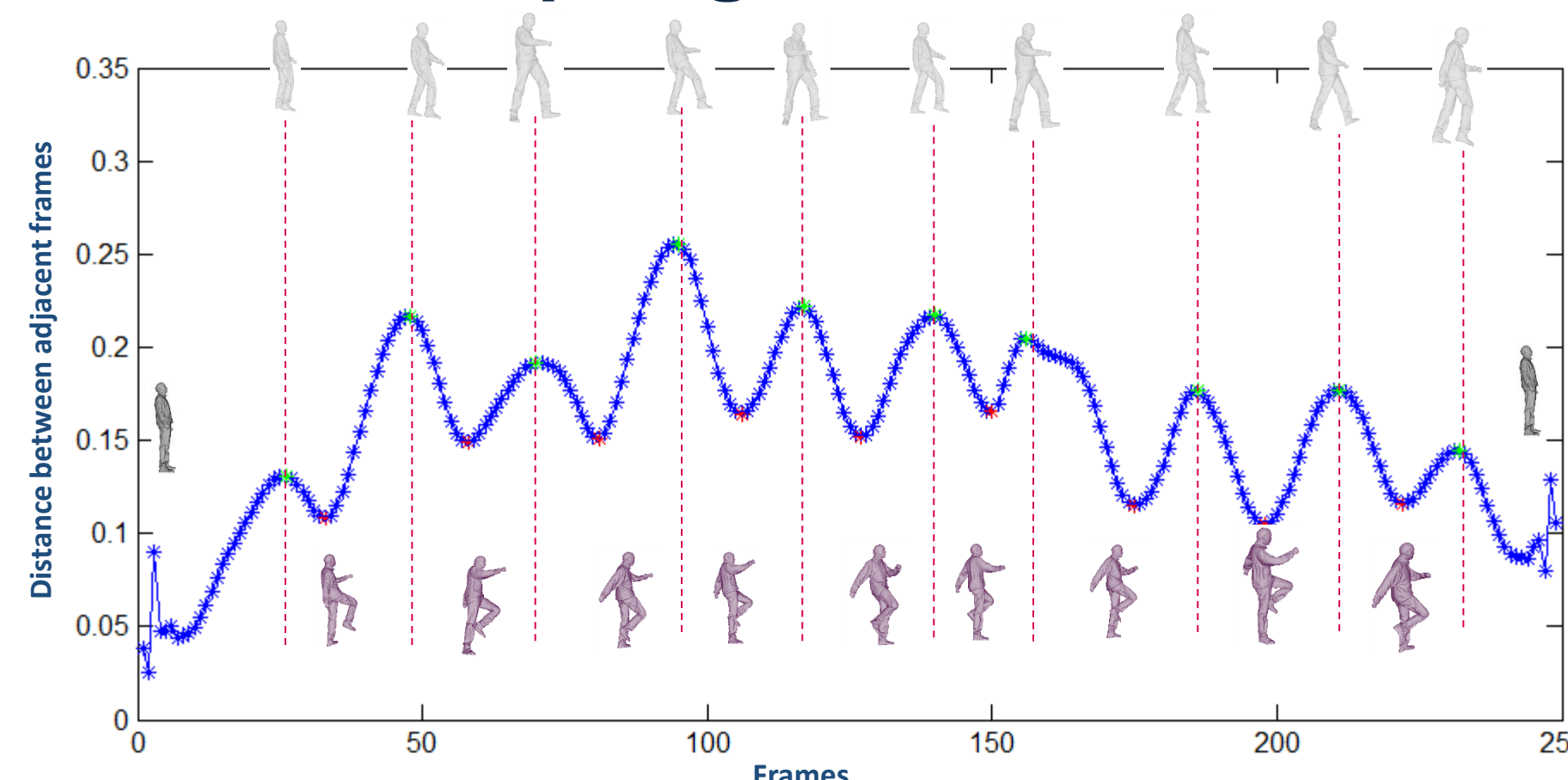


Results and Conclusion

Selected curves



Result of clip segmentation : Walk



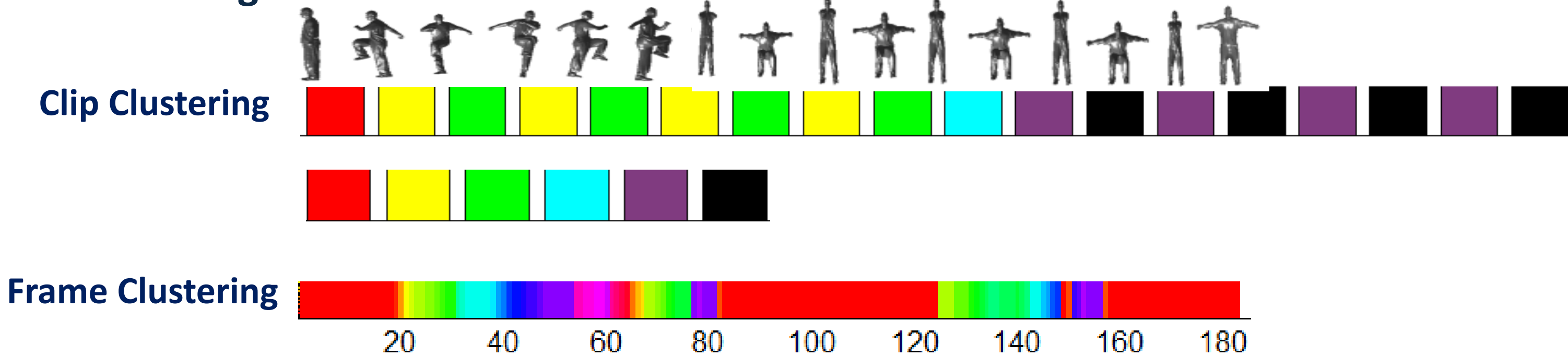
Conclusion

- Five curves are sufficient to represent at best the body pose
- Velocity curve is used to segment the long sequences into clips
- Clip matching using DTW on Riemannian manifold gives 93.44% of second tier rate and allows being invariant to speed
- Summarization by clustering is exploited in content-based motion retrieval

Selected publications

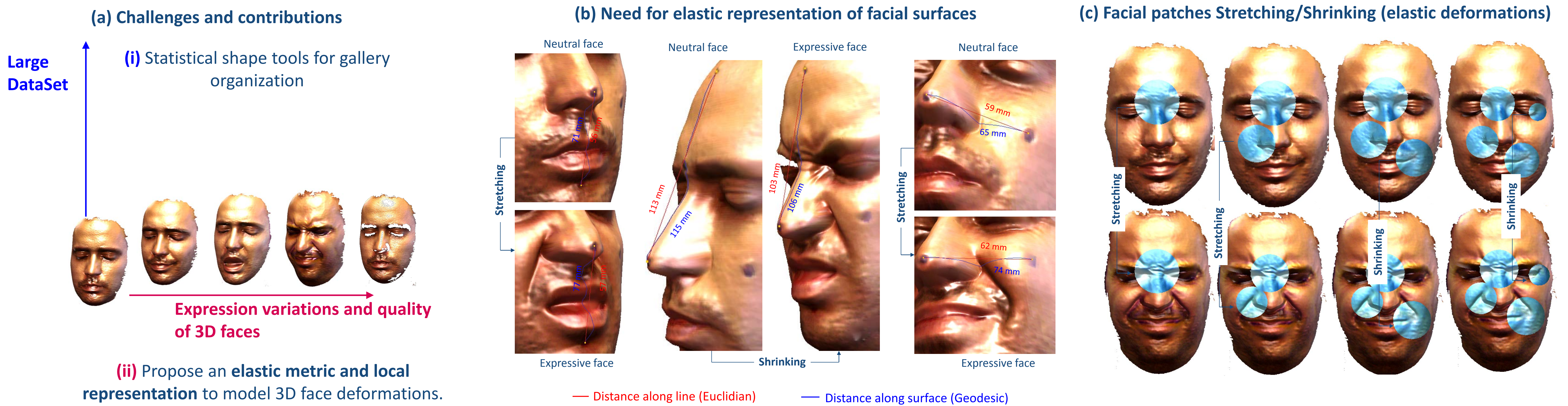
- R. Slama, H. Wannous, M. Daoudi, 3D human motion analysis framework for shape similarity and retrieval. Image Vision Computing Journal 32(2): 131-154 (2014)

Videoclustering



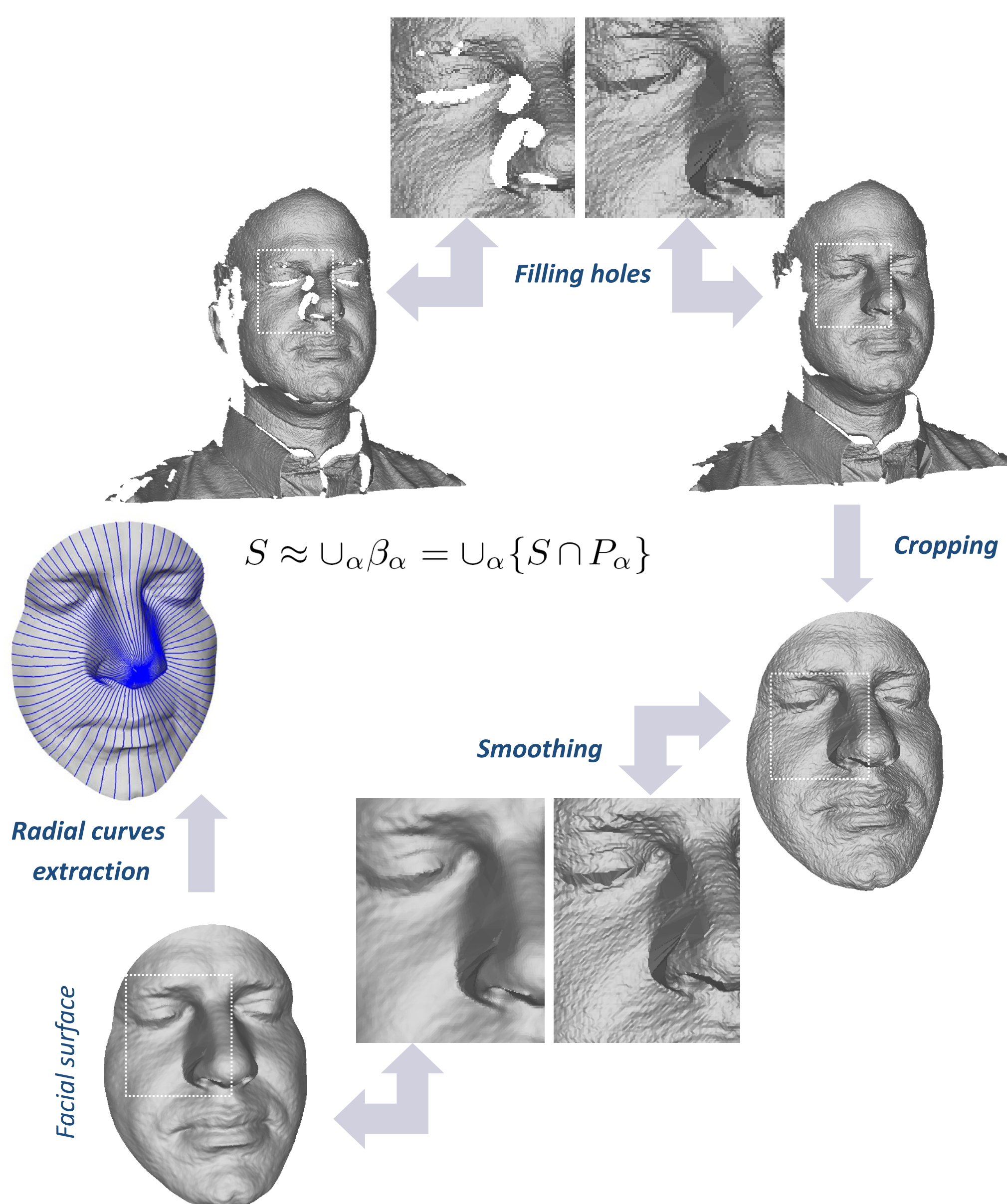
Statistical Shape Analysis of a Large 3D Faces Dataset

MOTIVATIONS

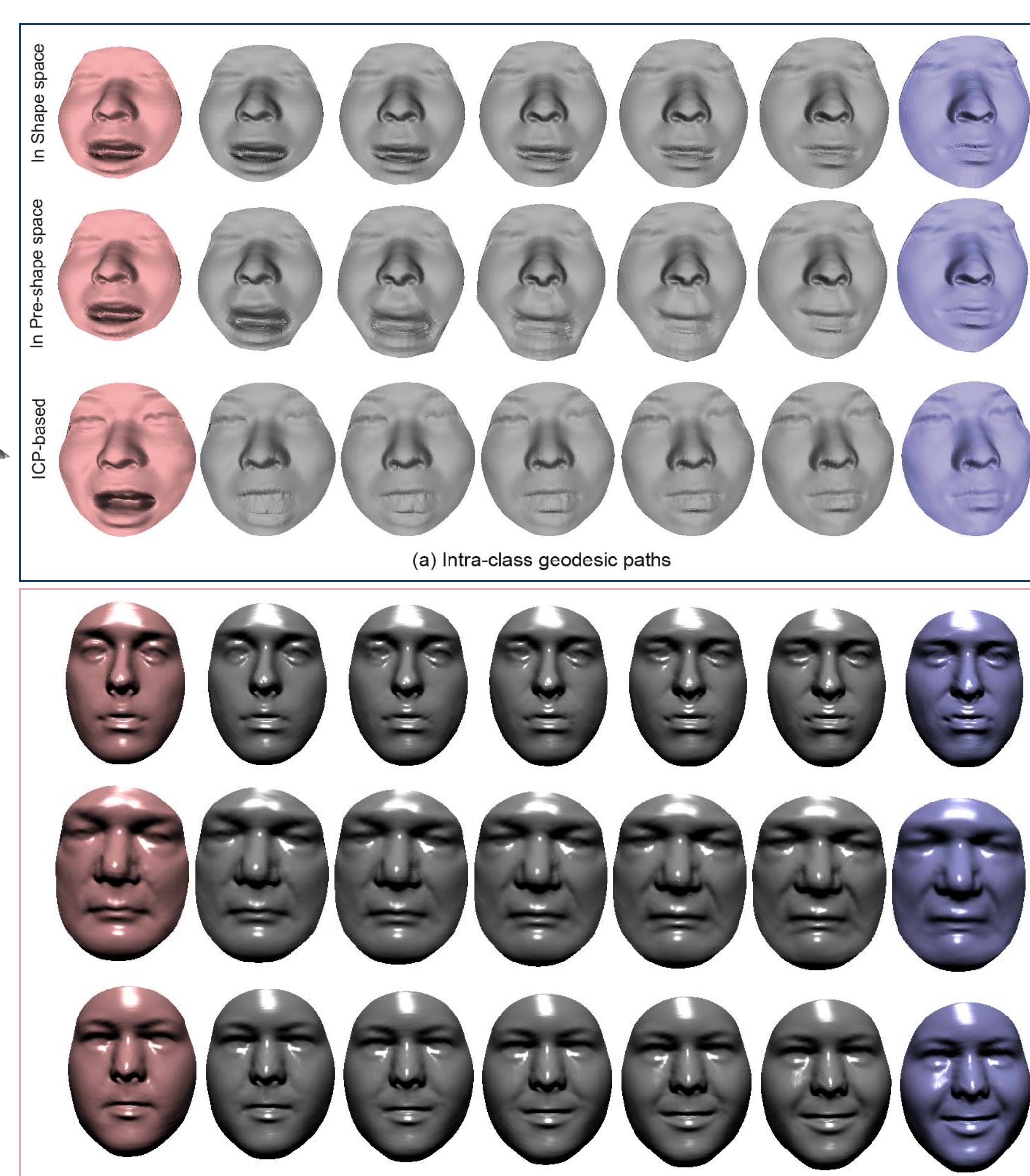


STATISTICAL SHAPE ANALYSIS of 3D FACE

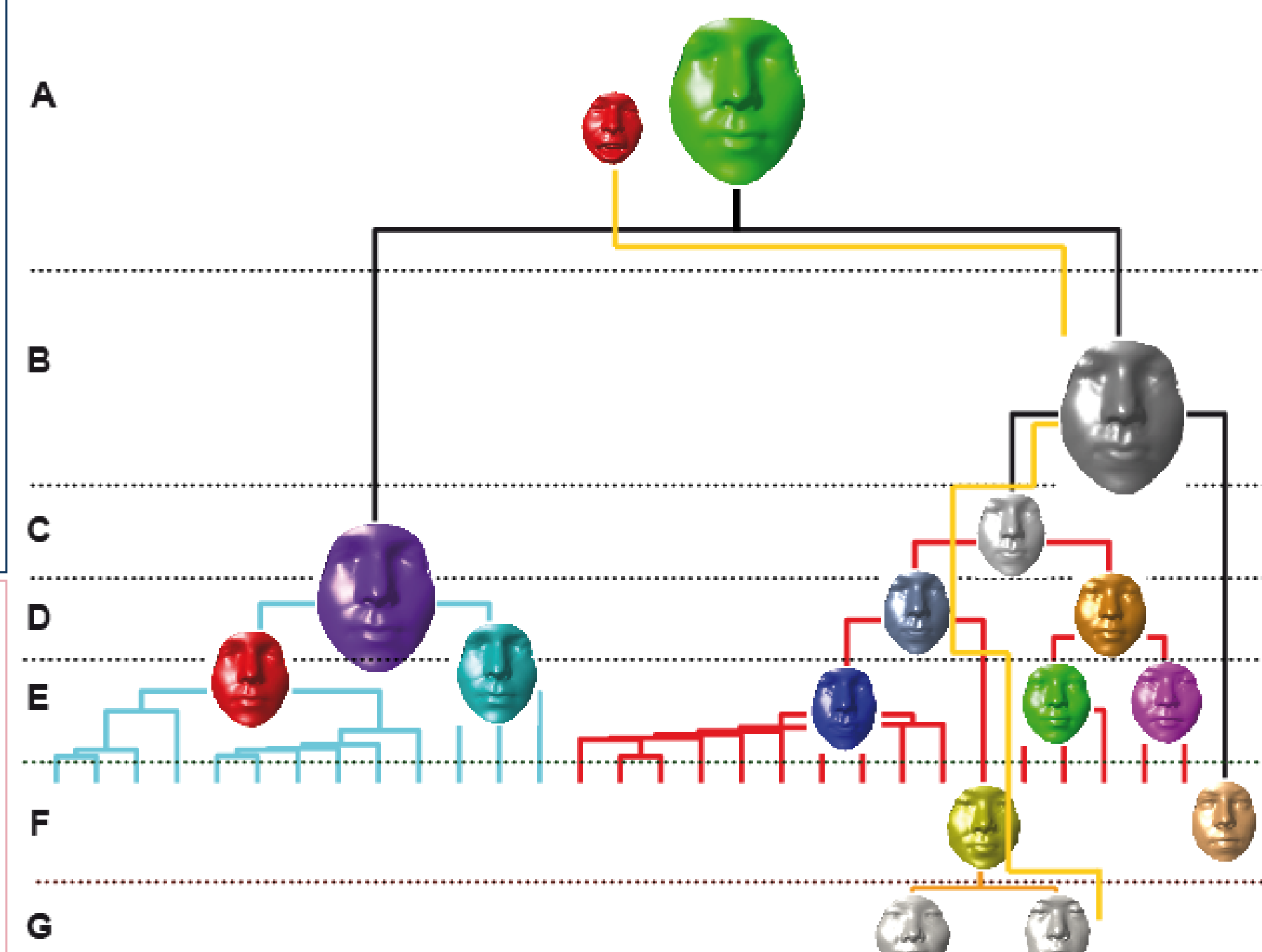
(a) Preprocessing of 3D scans and facial representation



(b) Examples of geodesics



(c) Gallery organization: Hierarchical retrieval of 3D faces Gallery



EXPERIMENTAL RESULTS

(a) Experimental protocol

- FRGCv2 dataset: 4007 facial scans
- 466 subjects
- Gallery: neutral faces of 466 subjects
- Probe: remaining

(b) Experimental results

Comparison of rank-1 scores on the FRGCv2 dataset with the state-of-the-art results.

Spreeuwiers [29]	Wang et al. [32]	Haar et al. [31]	Berretti et al. [2]	Queirolo et al. [26]	Faltemier et al. [9]	Kakadiaris et al. [13]	Our approach
99%	98.3%	97%	94.1%	98.4%	97.2%	97%	97%

Comparison of verification rates at FAR=0.1% on the FRGCv2 dataset with state-of-the-art results (the ROC III mask and the All vs. All scenario).

Approaches	Kakadiaris et al. [13]	Faltemier et al. [9]	Berretti et al. [2]	Queirolo et al. [26]	Spreeuwiers [29]	Wang et al. [32]	Our approach
ROC III	97%	94.8%	-	96.6%	94.6%	98.4%	97.14%
All vs. All	-	93.2%	81.2%	96.5%	94.6%	98.13%	93.96%

Selected publications

1. Hassen Drira, Boulbaba Ben Amor, Anuj Srivastava, Mohamed Daoudi, Rim Slama: 3D Face Recognition under Expressions, Occlusions, and Pose Variations. IEEE Trans. Pattern Anal. Mach. Intell. 35(9): 2270-2283 (2013).

2. Hassen Drira, Boulbaba Ben Amor, Anuj Srivastava, Mohamed Daoudi: A Riemannian analysis of 3D nose shapes for partial human biometrics. ICCV 2009: 2050-2057

IN A NUTSHELL

What is the performance of Bayesian bandit algorithms from a frequentist point of view? Bayes-UCB and Thompson Sampling appear to outperform frequentist algorithms on their own ground, which is supported by optimal regret bound for the Bernoulli case.

BAYESIAN VS. FREQUENTIST MODEL FOR MAB

K independent arms. Arm a depends on parameter θ_a and has expectation μ_a ; optimal arm is $a^* = \operatorname{argmax} \mu_a$ and $\mu^* = \mu_{a^*}$ is the highest expectation of reward associated.

Two probabilistic modelings

Frequentist :

- $\theta_1, \dots, \theta_K$ unknown parameters
- $(Y_{a,t})_t$ is i.i.d. with distribution ν_{θ_a}

Bayesian :

- $\theta_a \stackrel{i.i.d.}{\sim} \pi_a$
- $(Y_{a,t})_t$ is i.i.d. conditionally to θ_a with distribution ν_{θ_a}

At time t , arm A_t is chosen and reward $X_t = Y_{A_t,t}$ is observed

Two measures of performance

- Minimize (classic) regret

$$R_n(\theta) = \mathbb{E}_\theta \left[\sum_{t=1}^n \theta^* - \theta_{A_t} \right]$$

- Minimize “Bayesian“ regret

$$R_n = \int R_n(\theta) d\pi(\theta)$$

Optimal algorithms

- Asymptotically optimal algorithms satisfy, for $a : \mu_a < \mu^*$,

$$\limsup_{n \rightarrow \infty} \frac{\mathbb{E}_\theta [N_a(n)]}{\log(n)} \leq \frac{1}{\text{KL}(\nu_{\theta_a}, \nu_{\theta^*})}$$

They are optimal in the sense of Lai and Robbins' lower bound (1985) on the number of draw of a sub-optimal arm

- An index policy inspired by that of Gittins (1979) adapted to non-discounted rewards minimizes Bayesian regret

→ no frequentist guarantees

⇒ Our goal: Design algorithms inspired by the Bayesian modeling that are asymptotically optimal in the frequentist setting.

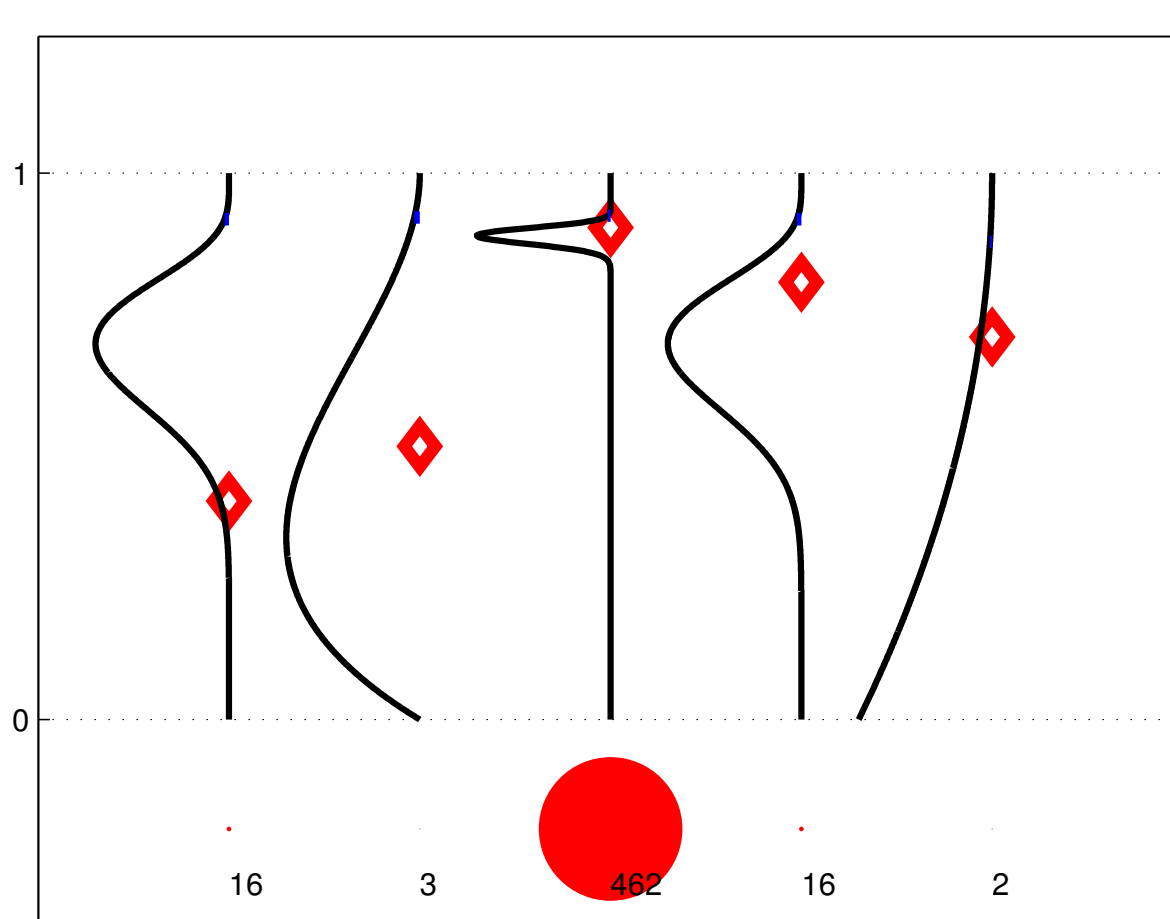
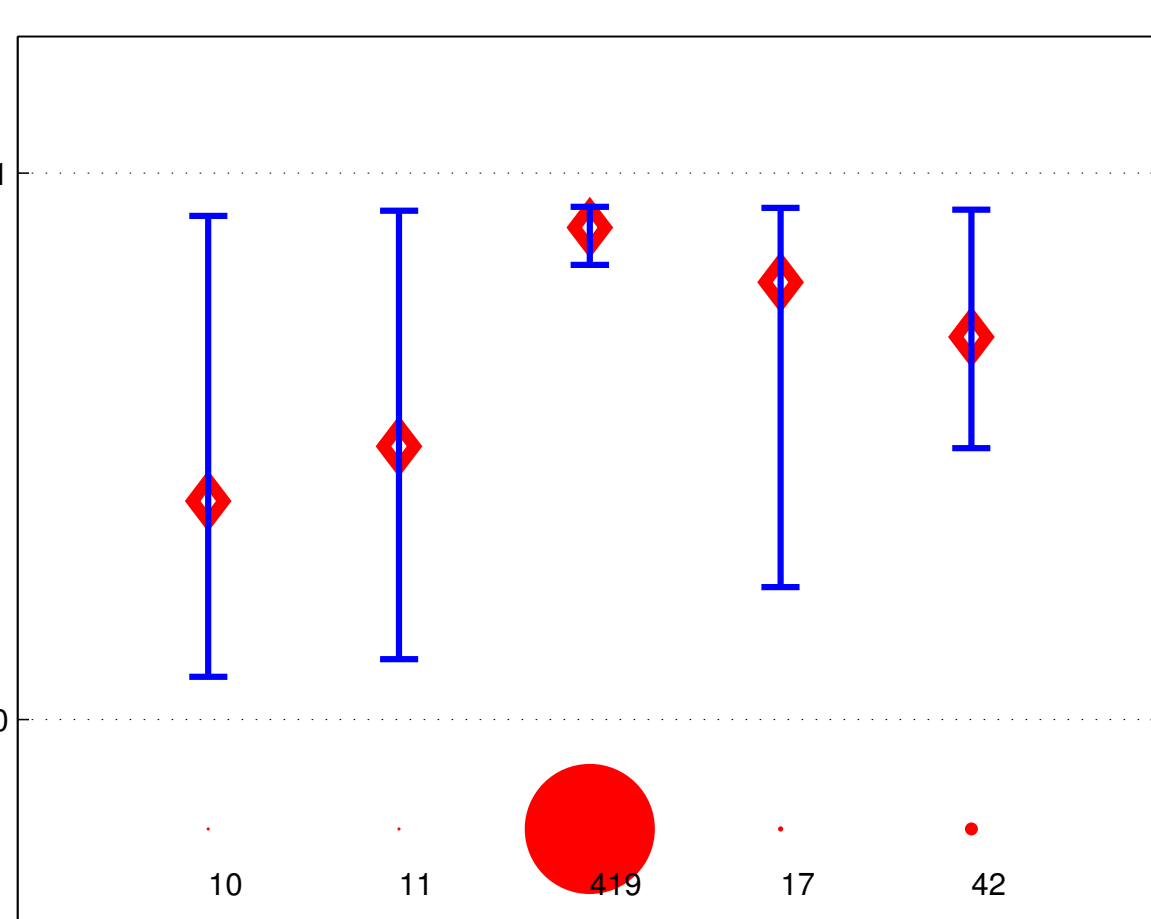
BAYESIAN VERSUS FREQUENTIST ALGORITHMS

Some quantities that naturally arise in the Bayesian modeling are

- $\Pi_t = (\pi_1^t, \dots, \pi_K^t)$ the current posterior over $(\theta_1, \dots, \theta_K)$
- $\Lambda_t = (\lambda_1^t, \dots, \lambda_K^t)$ the current posterior over the means (μ_1, \dots, μ_K)

Successful algorithms inspired by the frequentist modeling use

- Upper Confidence Bound for the empirical mean... (UCB)
- ... built using KL-divergence (KL-UCB, asymptotically optimal)



Whereas a Bayesian algorithm uses Π_t to determine action A_t .

BAYES-UCB AND THOMPSON SAMPLING

Bayes-UCB algorithm chooses $A_t = \operatorname{argmax}_{a=1..K} q_a(t)$, with

$$q_a(t) = Q \left(1 - \frac{1}{t(\log t)^c}, \lambda_a^t \right)$$

Thompson Sampling is a randomized algorithm:

$$\begin{cases} \forall a \in \{1..K\}, \theta_a(t) \sim \lambda_a^t \\ A_t = \operatorname{argmax}_a \theta_a(t) \end{cases}$$

Parameters: c (in practice, take $c = 0$), initial prior Π_0

BAYES-UCB: THEORETICAL ELEMENTS

ν_{θ_a} is the Bernoulli distribution $\mathcal{B}(\theta_a)$, π_a^0 the (conjugate) prior $\mathcal{U}([0, 1])$

- Bayes-UCB is asymptotically optimal for Bernoulli bandits

Theorem 1 Let $\epsilon > 0$; for the Bayes-UCB algorithm with parameter $c \geq 5$, the number of draws of a sub-optimal arm a is such that :

$$\mathbb{E}_\theta [N_a(n)] \leq \frac{1 + \epsilon}{\text{KL}(\mathcal{B}(\theta_a), \mathcal{B}(\theta^*))} \log(n) + o_{\epsilon,c}(\log(n))$$

- Bayes-UCB is very close to a frequentist algorithm

The Bayes-UCB index $q_a(t)$ is closely related to the one used by the KL-UCB algorithm (Cappé et al. 2013): $\tilde{u}_j(t) \leq q_j(t) \leq u_j(t)$ with:

$$u_a(t) = \operatorname{argmax}_{x > \frac{S_a(t)}{N_a(t)}} \left\{ d \left(\frac{S_a(t)}{N_a(t)}, x \right) \leq \frac{\log(t) + c \log(\log(n))}{N_a(t)} \right\}$$

$$\tilde{u}_a(t) = \operatorname{argmax}_{x > \frac{S_a(t)}{N_a(t)+1}} \left\{ d \left(\frac{S_a(t)}{N_a(t)+1}, x \right) \leq \frac{\log \left(\frac{t}{N_a(t)+2} \right) + c \log(\log(n))}{(N_a(t)+1)} \right\}$$

where $d(x, y) = \text{KL}(\mathcal{B}(x), \mathcal{B}(y)) = x \log \frac{x}{y} + (1-x) \log \frac{1-x}{1-y}$

Bayes-UCB appears to build automatically confidence intervals based on Kullback-Leibler divergence, that are adapted to the geometry of the problem in this specific case.

THOMPSON SAMPLING: THEORETICAL ELEMENTS

- TS is asymptotically optimal for Bernoulli bandits

Theorem 2 Let $\epsilon > 0$. With b defined below, for every sub-optimal arm a , there exists a constant $N(b, \epsilon, \theta_a, \theta^*)$ such that

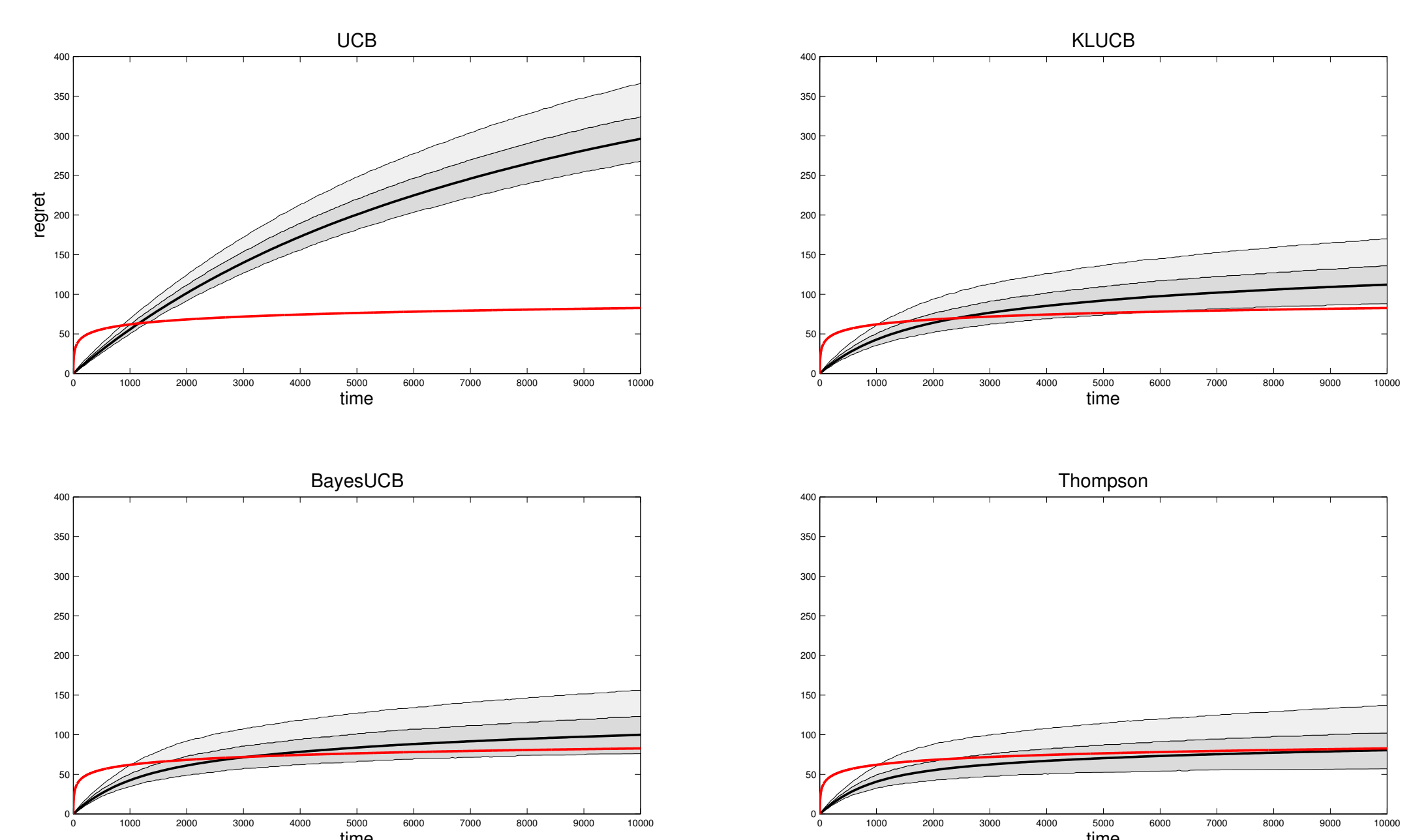
$$\mathbb{E}[N_a(n)] \leq (1 + \epsilon) \frac{\ln(n) + \ln \ln(n)}{d(\theta_a, \theta^*)} + N(b, \epsilon, \theta_a, \theta^*) + 5 + 2C_b.$$

Proof Bottleneck: For some constants $b = b(\mu) \in (0, 1)$ and $C_b < \infty$,

$$\sum_{t=1}^{\infty} \mathbb{P}(N_1(t) \leq t^b) \leq C_b.$$

NUMERICAL EXPERIMENTS AND BEYOND

- Bayesian algorithms are practically as efficient as optimal frequentist algorithms or even better!



Regret of the various algorithms as a function of time. The red curve show the lower bound, the black bold curve the mean regret and the dark and light shaded the central 99% and the upper 0.05%

- They are easier to implement: KL-UCB solves an optimization problem whereas Thompson Sampling only produces one sample!
- They are easy to generalize: general models where sampling from a posterior distribution is possible (using MCMC), sparse linear bandit, contextual bandit model...

REFERENCES

- [1] O. Cappé, A. Garivier, O. Maillard, R. Munos, G. Stoltz Kullback-leibler Upper Confidence Bounds for optimal sequential allocation *To appear in Annals of Statistics*, 2013
- [2] E. Kaufmann, O. Cappé, A. Garivier, Bayesian Upper Confidence Bounds for bandit problems *AISTATS*, 2012
- [3] E. Kaufmann, N. Korda, R. Munos Thompson Sampling: an asymptotically optimal finite-time analysis *ALT*, 2012



1. SUMMARY

The aggregation technique provides an estimator with well-established and excellent theoretical properties that applies for a wide family of times series which includes the $AR(d)$. However the numerical computation of this estimator relies on a Markov chain Monte Carlo method whose performances should be evaluated.

3. FORECASTERS

Let (X_1, \dots, X_n) observed values from this *stationary time series* $X = (X_t)_{t \in \mathbb{Z}}$. Consider a family of predictors $\{f_\theta, \theta \in \Theta\}$. For any $\theta \in \Theta$, f_θ is a function from which we obtain :

$$\hat{X}_t^\theta = f_\theta(X_{t-1}, \dots, X_{t-d}), \quad (2)$$

a possible forecasting of X_t according to θ .

Let ℓ be a loss function; we define the prediction risk as

$$R(\theta) = \mathbb{E} \left[\ell \left(\hat{X}_t^\theta, X_t \right) \right] \text{ and the empirical version of the risk as } r_n(\theta; X_1, \dots, X_n) = \frac{1}{n-d} \sum_{t=d+1}^n \ell \left(\hat{X}_t^\theta, X_t \right).$$

5. A THEORETICAL RESULT

Theorem 1 In the context of the $AR(d)$, for a bounded $\Theta \subset \mathbb{R}^p$, a uniform prior π yields that \exists a constant $\mathcal{E} : \forall \epsilon > 0$, with probability at least $1 - \epsilon$,

$$R \left(\hat{\theta}_{\sqrt{n}, n} \right) \leq \inf_{\theta \in \Theta} R(\theta) + \mathcal{E} \frac{\log^2(n)}{\sqrt{n}} + \frac{2}{\sqrt{n}} \log \left(\frac{1}{\epsilon} \right). \quad (4)$$

7. APPLICATION TO THE $AR(d)$ PROCESS

Since $s_d(1) \subseteq B_d(2^d - 1)$, the prior π can be defined on $\Theta = s_d(1)$ or $B_d(2^d - 1)$. These two possibilities are combined with two different proposals in the Metropolis-Hasting algorithm.

• Uniform prior on $B_d(2^d - 1)$

– *Uniform proposal* : $\beta_{\lambda, n} = \exp \left(-\lambda \mathcal{B}^2 \left(1 + \sqrt{d} (2^d - 1) \right)^2 \right)$.

– *Constrained random walk with Gaussian increment* :

$$\beta_{\lambda, n} = \left(\frac{n}{2\pi} \right)^{\frac{d}{2}} \exp \left(-2 (2^d - 1) \left(\lambda 2^{d+1} \mathcal{B}^2 + (2^d - 1) n \right) \right).$$

• Pushforward measure on $s_d(1)$

A map from the reciprocal roots of $\theta(z)$ into the coefficients $\theta_1, \dots, \theta_d$ and a measure on the first ones allow to define a prior on $s_d(1)$.

– *Uniform proposal* : $\beta_{\lambda, n} = \exp \left(-\lambda \mathcal{B}^2 \left(1 + \sqrt{d} (2^d - 1) \right)^2 \right)$.

– *Constrained random walk with Gaussian increment* :

$$\beta_{\lambda, n} = \left(\frac{n}{2\pi} \right)^{1+p+2\lfloor \frac{d}{2} \rfloor} \exp \left(-2 \left(\lambda 2^{d+1} (2^d - 1) \mathcal{B}^2 + \left(2^{1+d+2\lfloor \frac{d}{2} \rfloor} - 1 \right)^2 n \right) \right).$$

Theorem 3 \exists a constant \mathcal{F} such that $\forall m \geq M \left(\frac{\log(n)}{n}, \beta_{\sqrt{n}, n}, \epsilon \right)$, with M defined as in Theorem 2, with probability at least $(1 - \epsilon)^2$,

$$R \left(\bar{\theta}_{\sqrt{n}, n, m} \right) \leq \inf_{\theta \in \Theta} R(\theta) + \mathcal{F} \frac{\log^2(n)}{\sqrt{n}} + \frac{2}{\sqrt{n}} \log \left(\frac{1}{\epsilon} \right). \quad (7)$$

This result remains true for the whole family of Causal Bernoulli Shifts (CBS) processes (see [1] for the definition).

9. REFERENCES

References

- [1] Pierre Alquier and Olivier Wintenberger. *Model selection for weakly dependent time series forecasting*. Bernoulli, 18(3) : 883-913, 2012.
- [2] Krzysztof Łatuszyński and Wojciech Niemiro. *Rigorous confidence bounds for MCMC under a geometric drift condition*. J. Complexity, 27(1) : 23-38, 2011.

2. STABLE $AR(d)$ PROCESS

The $AR(d)$ is the stationary solution of :

$$X_t = \sum_{j=1}^d \theta_j X_{t-j} + \sigma \xi_t, \quad (1)$$

where the innovations ξ_t are i.i.d. with $\mathbb{E}\xi_t = 0$.

Denote by $s_d(1) = \left\{ (\theta_1, \dots, \theta_d) : \theta(z) = 1 - \sum_k \theta_k z^k \neq 0 \text{ for } |z| < 1 \right\}$.

We assume that the parameter $\bar{\theta} = (\theta_1, \dots, \theta_d) \in s_d(1)$ and that (ξ_t) have compact support $\Rightarrow \exists \mathcal{B} \in \mathbb{R}_+^* : X_t \in [-\mathcal{B}, \mathcal{B}] \forall t$.

4. GIBBS ESTIMATOR

For a $\lambda > 0$ (temperature parameter), we define the *Gibbs estimator* as the expectation of a r.v. drawn under the Gibbs measure $\pi \{-\lambda r_n\}$:

$$\hat{\theta}_{\lambda, n} = \pi \{-\lambda r_n\} [\text{Id}] = \int_{\Theta} \theta \pi \{-\lambda r_n(\cdot)\} (d\theta), \quad (3)$$

where $\nu[h] = \int h d\nu$ and $\nu\{h\}(d\theta) = \frac{\exp(h(\theta))}{\nu[\exp(h)]} \nu(d\theta)$.

6. NUMERICAL APPROXIMATION [2]

The *Metropolis-Hastings algorithm* generates a Markov chain $\Phi = \{\Phi_i\}_{i \geq 0}$ with the target distribution ρ as a unique invariant measure, based on another Markov chain which serves as a proposal. We tested :

- The *independent Hastings algorithm* where the proposal is i.i.d. with density q such that $\frac{q(y)}{\rho(y)} \geq \beta, \forall y \in \Theta$ for some $\beta > 0$.
- The *Metropolis-Hastings algorithm* where the proposal is a Markov chain with conditional density kernel q on $\bar{\Theta} \times \bar{\Theta}$ such that $\beta = \inf_{x \in \bar{\Theta}, y \in \bar{\Theta}} \frac{\rho(y)}{\rho(x)} \inf_{x \in \bar{\Theta}, y \in \bar{\Theta}} q(x, y) > 0$.

$\bar{\theta}_m = \frac{1}{m} \sum_{i=0}^{m-1} \Phi_i$ is a numerical estimate of $\int x \rho(x) dx$.

Theorem 2 Note by $\text{diam}(\Theta) = \sup_{x, y \in \Theta} \|x - y\|$ and define :

$$M(\alpha, \gamma, \epsilon) = \frac{(2 - \gamma) \text{diam}(\Theta)}{2\alpha^2 \epsilon \gamma} + \frac{1}{2} \sqrt{\left(\frac{(2 - \gamma) \text{diam}(\Theta)}{\alpha^2 \epsilon \gamma} \right)^2 + \frac{4 \text{diam}(\Theta)}{\alpha^2 \epsilon \gamma}}, \quad (5)$$

In the two previous cases, $\forall m \geq M(\alpha, \beta, \epsilon)$, with probability at least $1 - \epsilon$,

$$\left| \bar{\theta}_m - \int x \rho(x) dx \right| \leq \alpha. \quad (6)$$

8. NUMERICAL RESULTS

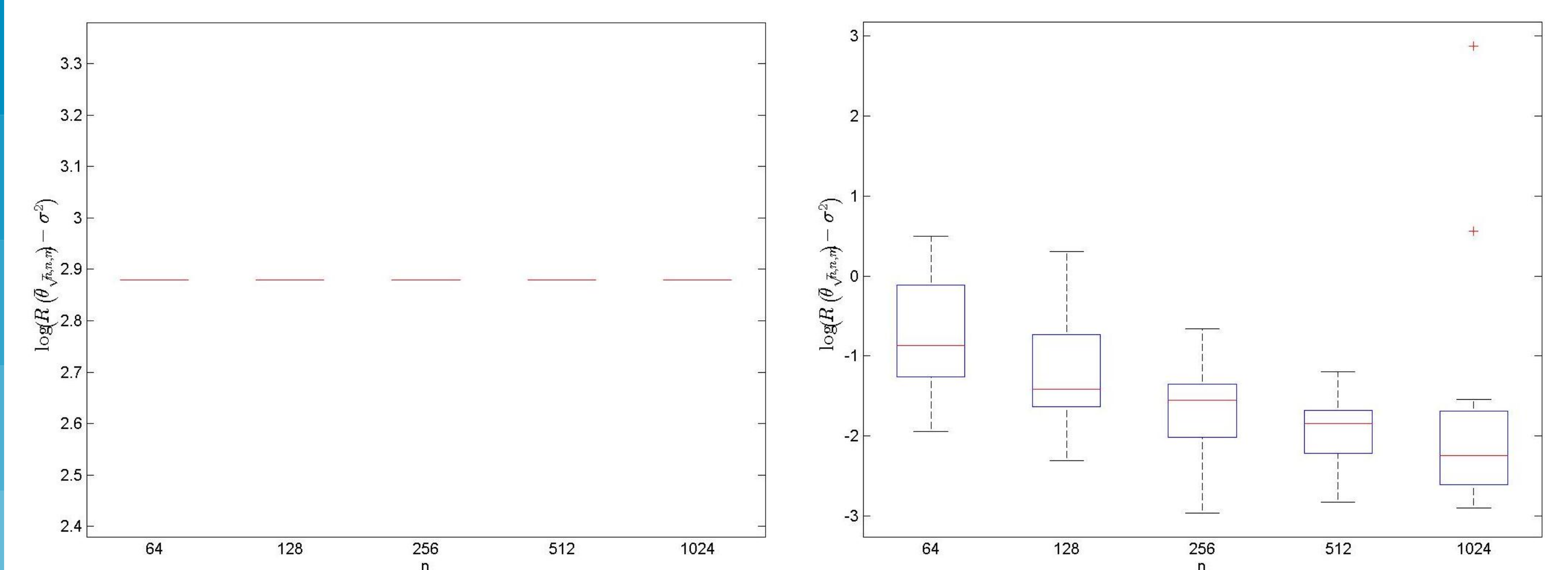


Figure 1: Uniform proposal, $d = 8, \Theta = B_8(2^8 - 1)$

Figure 2: Gaussian proposal, $d = 8, \Theta = s_8(1)$

Figure 2 shows good results in contrast to Figure 1. However, using (5) and the obtained expressions of β yield to the following equivalence for the minimal number of iterations m guarantying a correct prediction error :

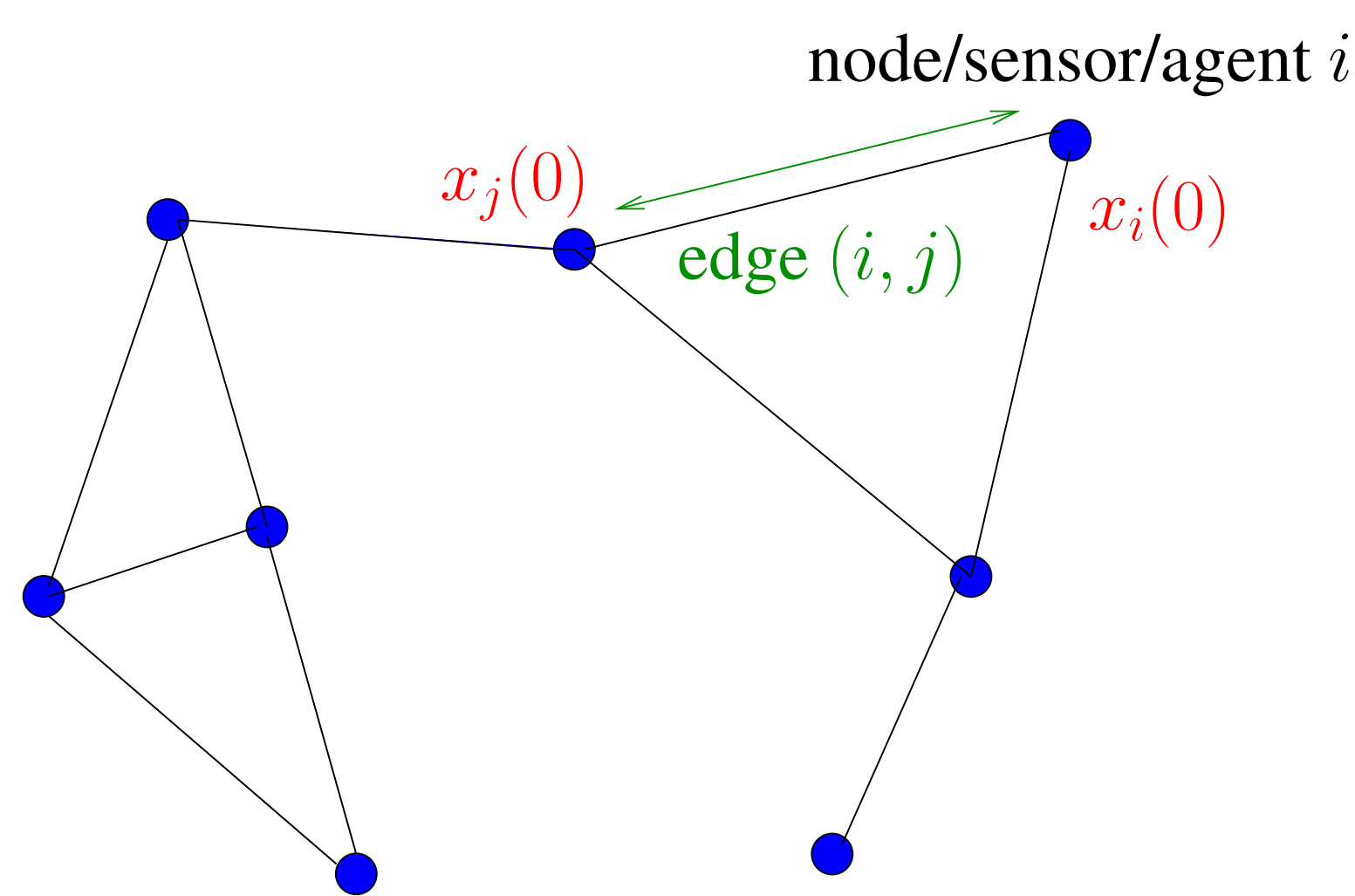
$$m \geq C_1(d) \frac{\log^2 n}{2n^2 \epsilon} \exp(C_2(d) \sqrt{n}),$$

where C_1 and C_2 are positive functions. Hence, to guaranty the Formula (7) we need a very high number of iterations, which can lead to a prohibitive computational cost.

ABSTRACT

We study the distributed computation of the average value of initial measurements in a Wireless Sensor Network. Unlike existing works, we take benefit of the broadcast nature of the wireless channel to speed up the convergence speed

MODEL



- N nodes/sensors/agents
- $x_i(0)$ measurement at node i at time 0
- **Problem:** at each node, we want to compute

$$x_{ave} = \frac{1}{N} \sum_{i=1}^N x_i(0)$$

without fusion center and so with only local communications.

- **Applications:** practical measurements of temperature, gas pressure,
- At each time t , the sensors' updates can be rewritten in matrix form

$$\mathbf{x}(t+1) = \mathbf{K}(t)\mathbf{x}(t)$$

where $\mathbf{K}(t)$ is the update matrix for time t

Goal:

- We want $\mathbf{x}(t)$ to converge to the average consensus $x_{ave}\mathbf{1}$
- We want to wake up only one node at each time
- We want to use the broadcast nature of the wireless links

AVERAGING ALGORITHMS

Fundamental properties to ensure convergence to the average consensus:

Property	Mathematical implication
1. Consensus preservation $\mathbf{x}(t) = c\mathbf{1} \Rightarrow \mathbf{x}(t+1) = c\mathbf{1}$	Row stochasticity $\mathbf{K}\mathbf{1} = \mathbf{1}$
2. Consensus value $\mathbf{x}(t) = c\mathbf{1} \Rightarrow \mathbf{x}(t) = x_{ave}\mathbf{1}$	Column stochasticity $\mathbf{1}^T\mathbf{K} = \mathbf{1}^T$
3. Convergence $\mathbb{E}[\ \mathbf{x}(t) - c\mathbf{1}\ ^2] \rightarrow 0$	Spectral radius $\rho(\mathbb{E}[\mathbf{K}] - 1/N\mathbf{1}\mathbf{1}^T) < 1$

Consequences:

- Double-stochasticity needs feedback, so **no broadcast**.
- If only column-stochasticity, $\exists \mathbf{v}(t)$ non-negative vectors s.t.

$$\mathbf{K}(t) \cdots \mathbf{K}(2)\mathbf{K}(1) \sim \mathbf{v}(t)\mathbf{1}^T \quad \text{and} \quad \mathbf{x}(t) \sim (N x_{ave})\mathbf{v}(t)$$

Proposed solution:

- An other variable updated with the same matrix has to be considered to know $\mathbf{v}(t)$. Then $\mathbf{v}(t)$ can be removed by division.

⇒ **Variables do not converge to consensus, but the quotient does.**

PROPOSED ALGORITHM: BWGossip

Algorithm based on the *Sum-Weight* framework [Kempe2003, Bénézit2011] where sensors have **two** local variables jointly updated:

- a sum variable $\mathbf{s}(t)$ $\begin{cases} \mathbf{s}(t+1) = \mathbf{K}(t)\mathbf{s}(t) \\ \mathbf{w}(t+1) = \mathbf{K}(t)\mathbf{w}(t) \end{cases}$ and $\mathbf{s}(0) = \mathbf{x}(0)$
- a weight variable $\mathbf{w}(t)$ $\mathbf{w}(0) = \mathbf{1}$

BWGossip

Assuming that at time t , the sensor i wakes up

- ▶ Sensor i broadcasts $\left(\frac{s_i(t)}{|\mathcal{N}_i|+1}, \frac{w_i(t)}{|\mathcal{N}_i|+1}\right)$
- ▶ At sensors in the neighborhood \mathcal{N}_i , we have:

$$\begin{cases} s_j(t+1) = s_j(t) + \frac{s_i(t)}{|\mathcal{N}_i|+1} \\ w_j(t+1) = w_j(t) + \frac{w_i(t)}{|\mathcal{N}_i|+1} \end{cases}, \forall j \in \mathcal{N}_i$$

- ▶ At sensor i , we have :

$$\begin{cases} s_i(t+1) = \frac{s_i(t)}{|\mathcal{N}_i|+1} \\ w_i(t+1) = \frac{w_i(t)}{|\mathcal{N}_i|+1} \end{cases}$$

- ▶ All other sensors stay idle.

RESULTS

Theorem 1 CONVERGENCE

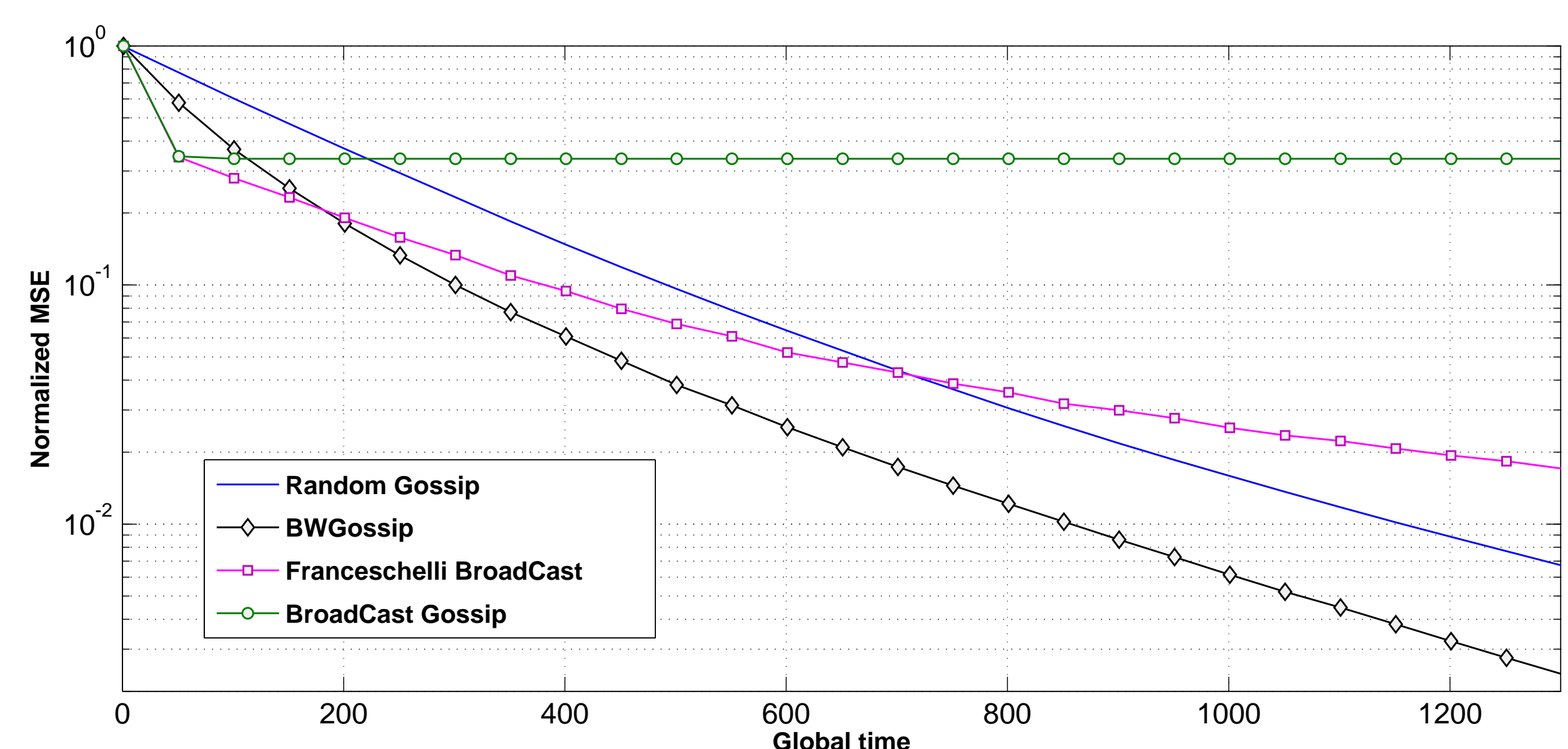
BWGossip converges to the average consensus almost surely.

Theorem 2 CONVERGENCE SPEED (main result)

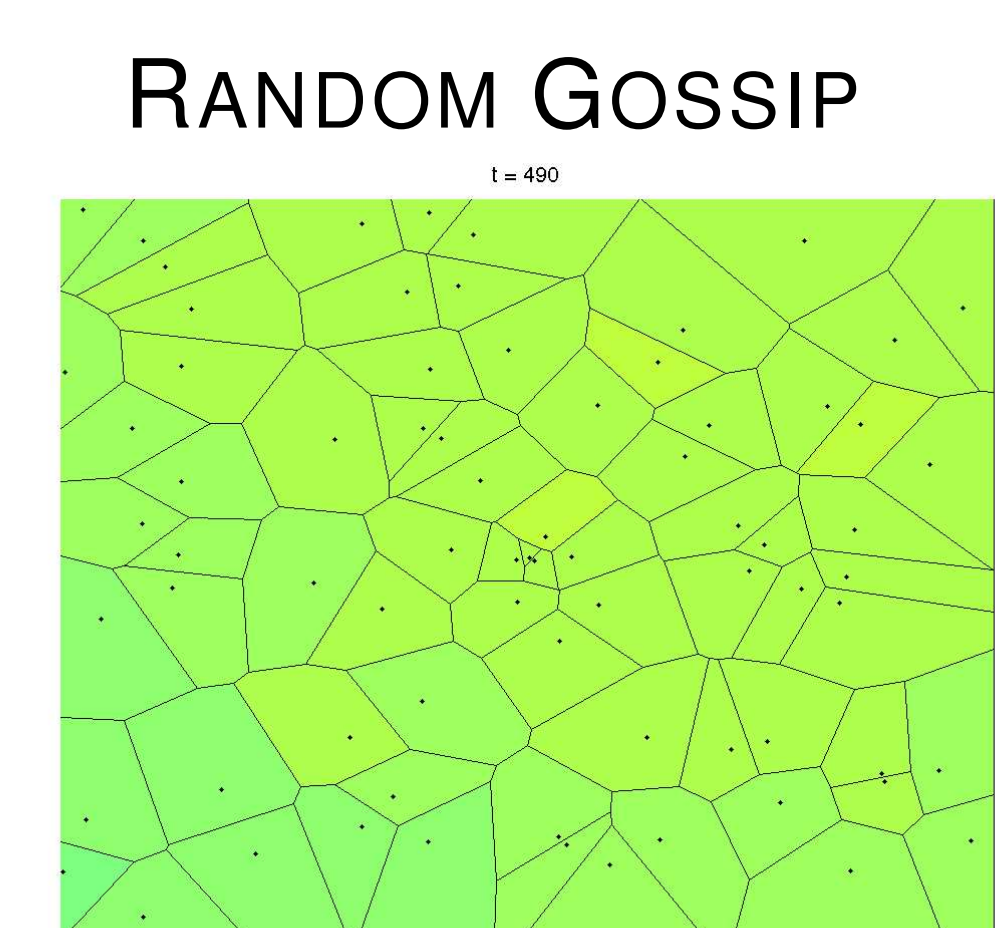
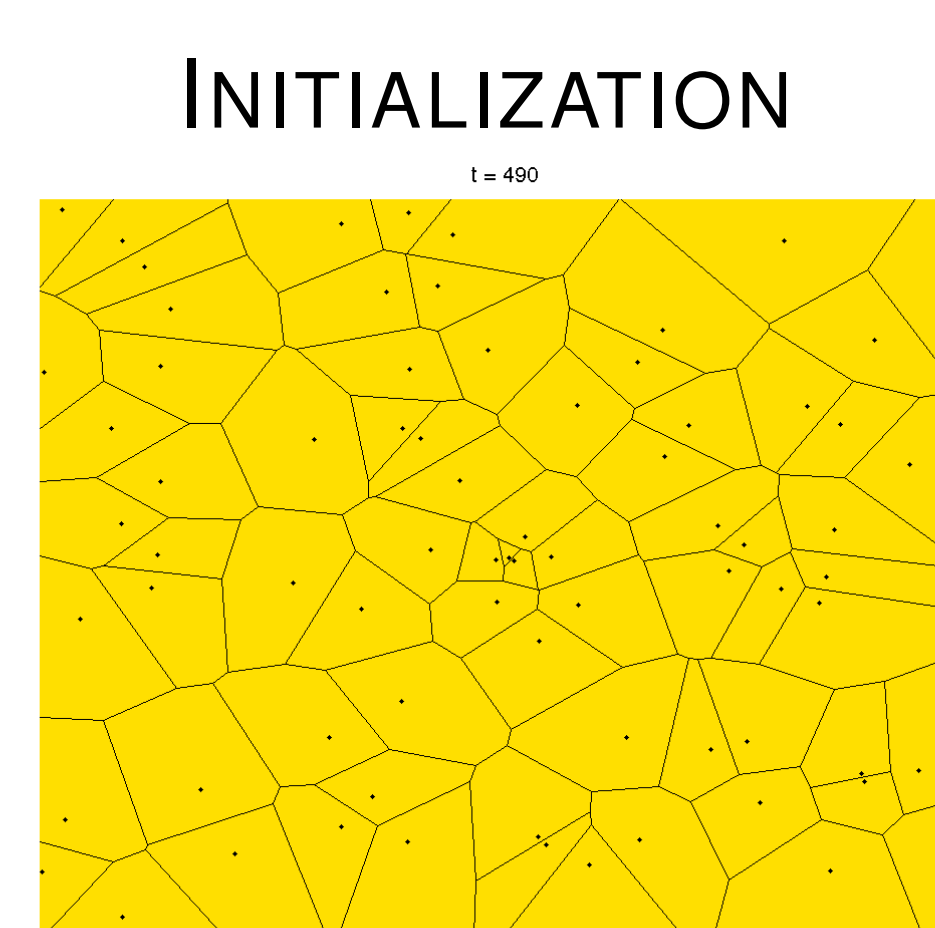
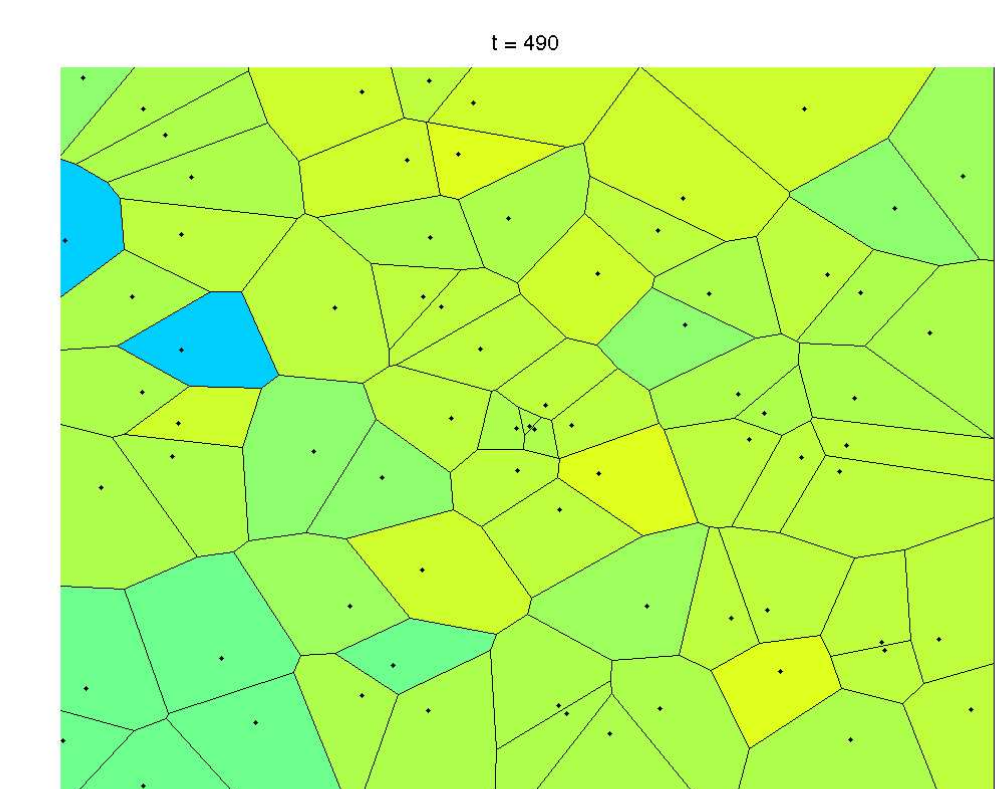
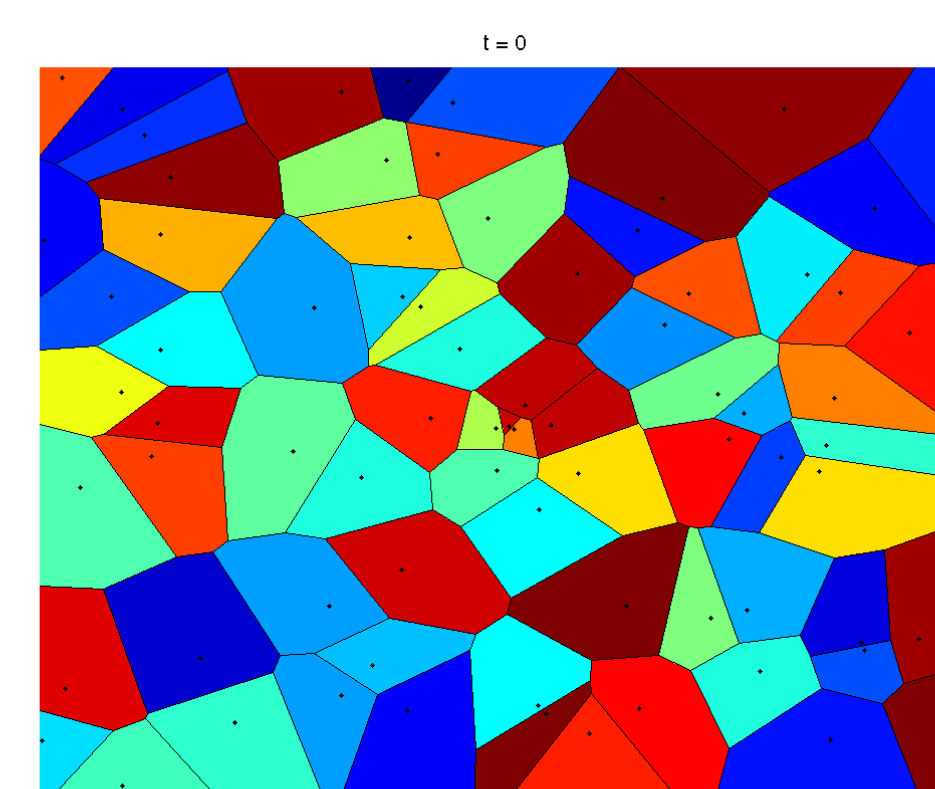
For the BWGossip algorithm, we have:

$$\forall \epsilon > 0, \quad \|\mathbf{x}(t) - x_{ave}\mathbf{1}\|_2^2 = \mathcal{O}_P((\Gamma + \epsilon)^t)$$

with $\Gamma = \rho((\mathbf{I} - \mathbf{J}) \otimes (\mathbf{I} - \mathbf{J}) \cdot \mathbb{E}[\mathbf{K} \otimes \mathbf{K}]) < 1$



- **BWGossip outperforms the existing algorithms.**



BROADCAST GOSSIP

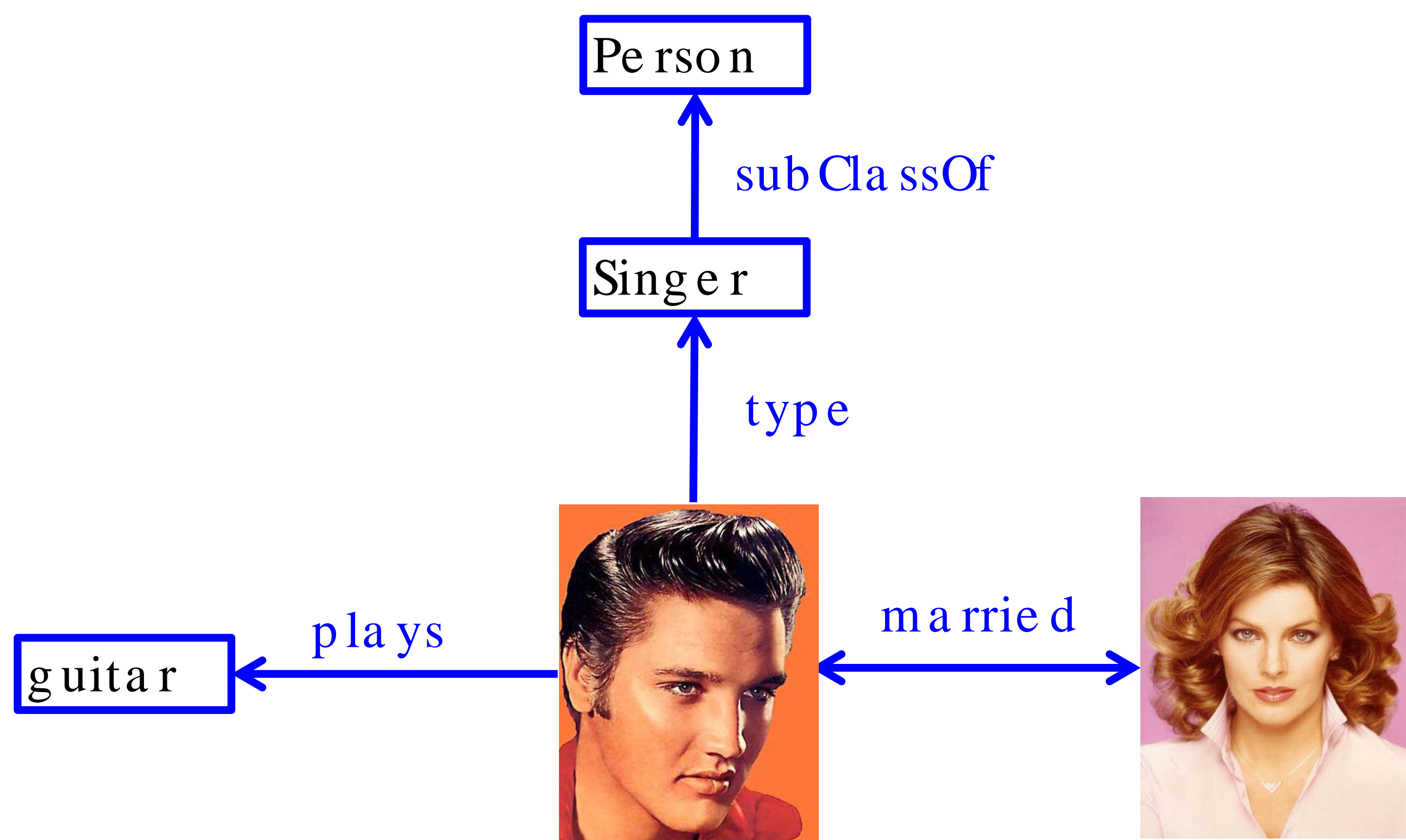
BW GOSSIP

The Independent Otto Hahn Research Group

Ontologies

Fabian M. Suchanek

Ontology Construction



Ontologies serve, e.g., for disambiguation, translation, and question answering. We develop YAGO, the largest public ontology with a quality guarantee. YAGO is built automatically from Wikipedia and other sources.
<http://yago-knowledge.org>

YAGO – A Core of Semantic Knowledge

Fabian M. Suchanek, Gjergji Kasneci, Gerhard Weikum (WWW 2007)
 + follow-up publications in 2008, 2011, 2012, 2013

Rule Mining



? \rightarrow $\text{popSinger}(x) \Rightarrow \text{is}(x, \text{rich})$

We develop techniques to mine rules, correlations, and schemas from an ontology. These serve to propose missing links, detect inconsistencies, reveal correlations, and make sense out of data. The semantics, incompleteness, and the scale of the data are different from classical settings.

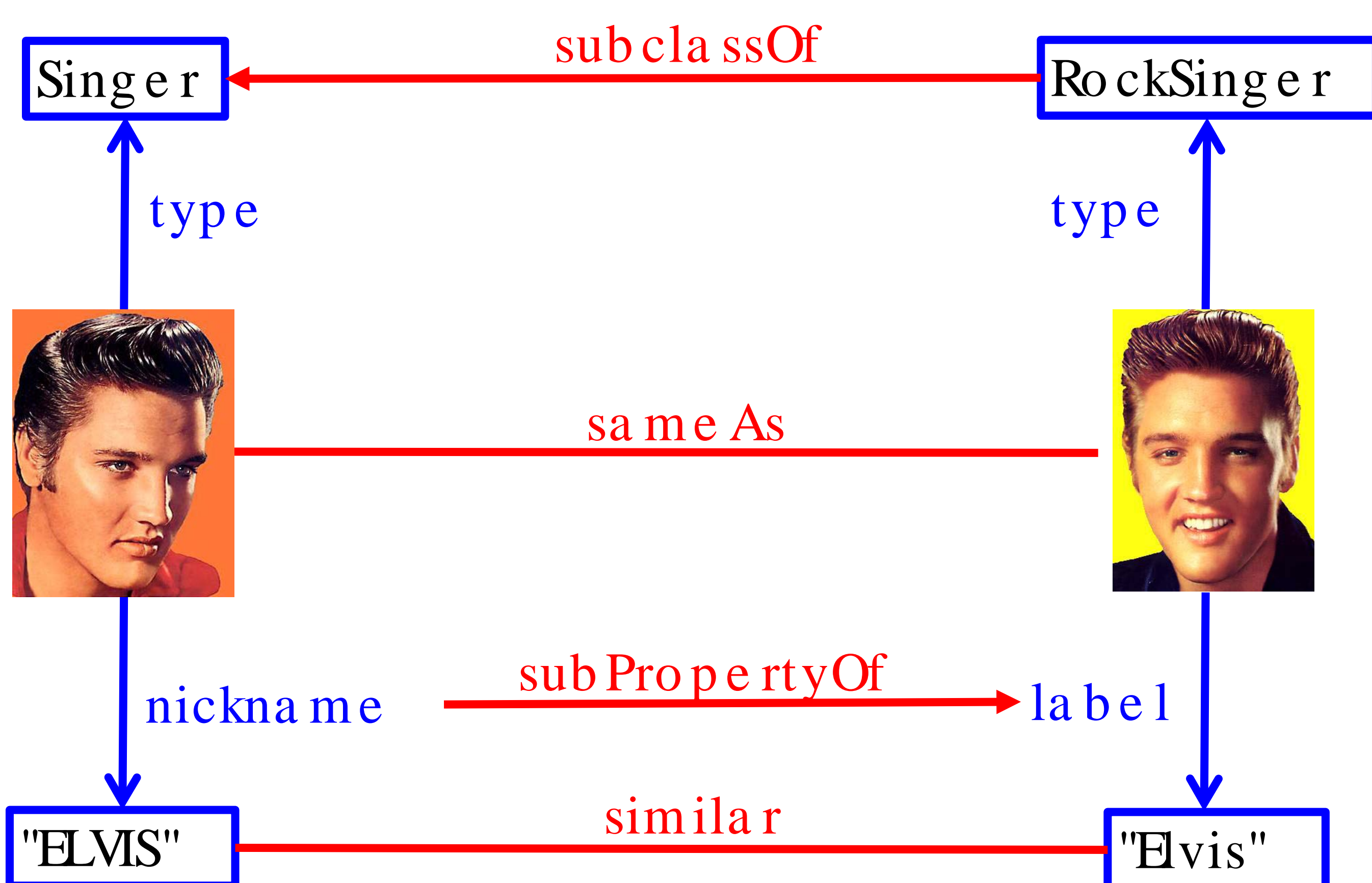
Some rules we mined on YAGO are:

$\text{hasChild}(x,y) \ \& \ \text{hasChild}(z,y) \Rightarrow \text{married}(x,z)$
 $\text{wonAward}(x, \text{LeibnizPreis}) \Rightarrow \text{livesIn}(x, \text{Germany})$
 $\text{acadAdvisor}(x,y) \ \& \ \text{almaMater}(y,z) \Rightarrow \text{worksAt}(x,z)$

AME Association Rule Mining under Incomplete Evidence

Luis Galárraga, C.Teflioudi, KHse, F.M.Suchanek (WWW 2013)

Ontology Matching

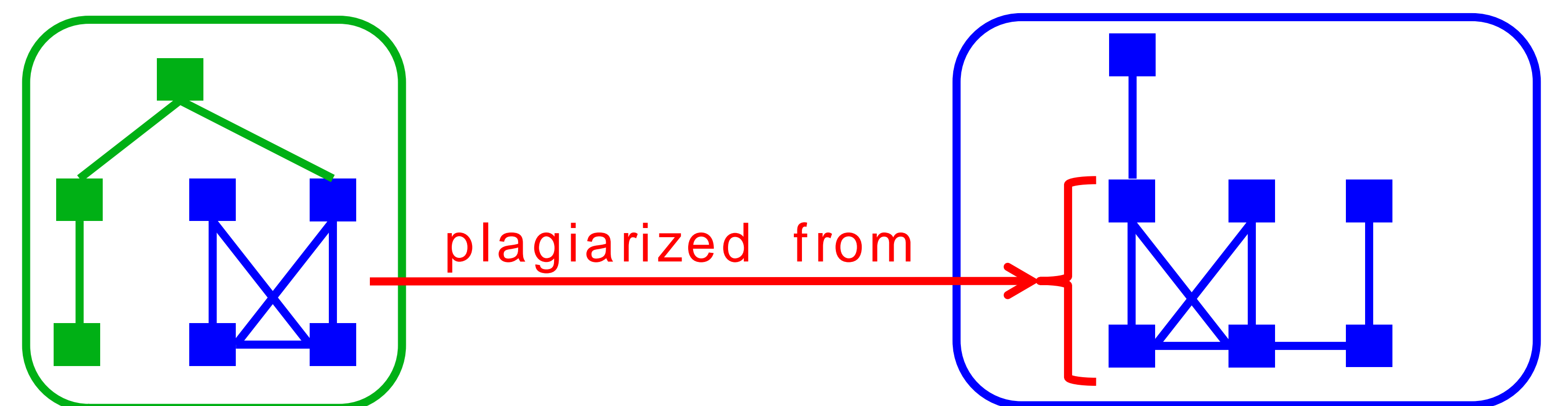


Many ontologies on the Semantic Web contain information about the same entities. To make use of complementary information, one has to determine which entities, classes, literals, and properties correspond. We develop statistical, logical, and probabilistic models and algorithms for this purpose.

PARIS Probabilistic Alignment of Relations, Instances and Schema

Fabian M. Suchanek, S. Abiteboul, P. Senellart (VLDB 2012)

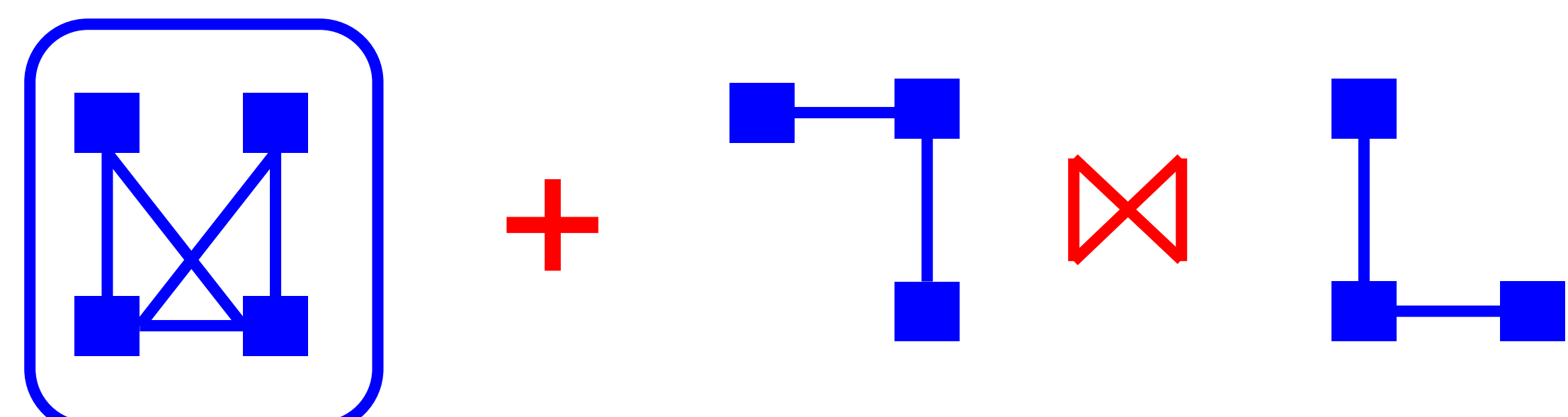
Provenance



Statements can be illegally copied over from another ontology. We develop techniques to trace plagiarisms.

Watermarking for Ontologies (ISAC 2011)

Fabian M. Suchanek, David Gross-Amblard, Serge Abiteboul



We develop logic-based models to integrate Web services into ontologies.

SUSE Search using Services and Information Extraction

Ncoleta Preda, F. M. Suchanek, W. Yuan, G. Weikum (ICDE 2013)



MAX-PLANCK-GESELLSCHAFT

Casting a Web of Trust over Wikipedia: an Interaction-based Approach

Silviu Maniu, Talel Abdesslem, Bogdan Cautis; Télécom ParisTech – CNRS LTCI, Paris, France, {firstname.lastname@telecom-paristech.fr}

Our goal

Uncover a **signed network** over Wikipedia contributors from their **interactions**.

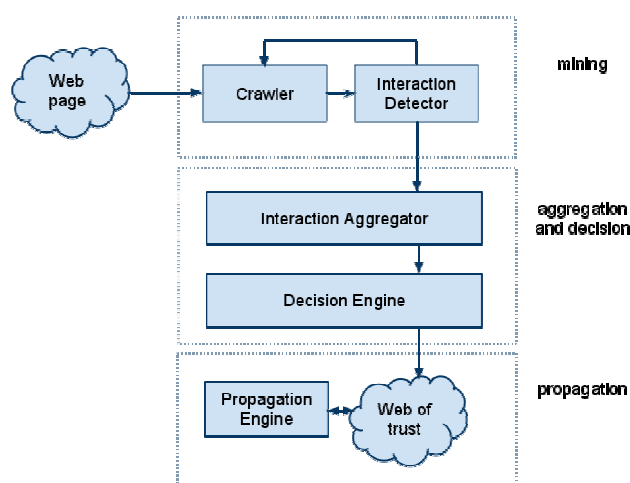
Signed link = a representation of the degree of trust/similarity or distrust/dissimilarity between two users

Several signed networks are already present in social media: Epinions (trust/distrust tags), Slashdot (friend/foe), Wikipedia Elections (support/oppose votes)

Motivation:

- Social applications can be enhanced by knowing such signed links (social search systems, recommender systems, trust and reputation, etc.)

General Architecture



Dataset

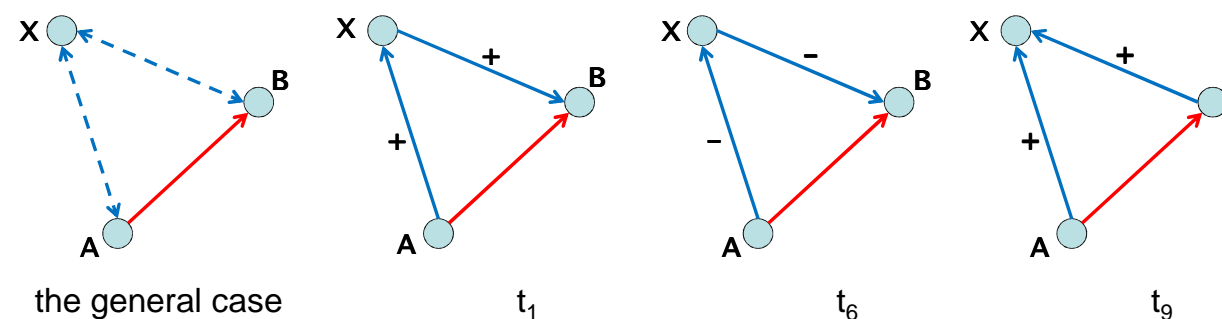
From the politics domain of the English Wikipedia:

- we extracted 320 articles, for a total of 442,297 revisions submitted by 105,177 contributors
- resulting in 800,057 total interactions, in which participate 42,631 adminship votes and 2,913 barnstars

Validation

Does this network represent a *plausible configuration of link signs*?

First approach: comparing our network with three existing explicit signed networks (Epinions, Slashdot, Wikipedia Elections)



Using the concept of **link triads** and the predictions of two social theories: balance and status, for measuring:

The global properties of WikiSigned

- Our network has similar triad distribution as the explicit networks
- And it has the global structure of a network in which status theory holds (only one contradiction for triad signs)

The local properties of WikiSigned

- It can self-predict its link signs with 0.822 accuracy (AUC of 0.899)
- Good accuracy also in cross training-predicting (training on the row data and predicting on the column data)

	Epinions	Slashdot	Elections	WikiSigned
Epinions	0.906	0.905	0.787	0.727
Slashdot	0.929	0.806	0.792	0.732
Elections	0.922	0.895	0.814	0.733
WikiSigned	0.889	0.844	0.784	0.822

Second approach: application-level validation

- Predicting the importance or quality of articles by using the **knowledge of link composition** (number of positive and negative links) in training predictive models
- This knowledge of link signs helps the prediction when we predict the article importance

feature	Quality	Importance
Contributors	0.683	0.566
Contributors + normal links	0.740	0.779
Social links	0.807	0.750

Interactions in Wikipedia

Article editing history for 'Paralympiakos' showing various actions like 'Reverted edit', 'Adding a parameter', 'Removing a parameter', etc.

Article editing:

- on text content (**inserting**, **deleting** and **replacing** text between the contributors)
- on the article revisions (**reverting**/discarding a version of an article and **restoring** another)

Wikipedia:Requests for adminship/Dianna. Shows a nomination process with comments, a 'Good job' nomination, and the awarding of 'The Working Man's Barnstar' and 'The Random Acts of Kindness Barnstar'.

Adminship election:

- Contributors participate in so called requests for adminship, elections in which contributors can:
 - Support** the candidate
 - Oppose** the candidate

Interactions on user pages:

- Contributors can give each other prizes called **barnstars** (generally for good behavior)

Deciding link signs



The **interaction vector** = aggregation of all interactions for each pair of users

The decision process: -1 (negative) or +1 (positive) link

- Annotate the atomic interactions with signs (positive or negative, as shown above)
- Each interaction votes with its weight (measure) for the sign of the higher-order type
- All types vote for the final link sign

WikiSigned - the resulting network

- 71,770 nodes and 463,312 edges, of which 85.93% positive

A Privacy Management System for Social Networks

Imen BEN DHIA, Talel ABDESSALEM, Mauro SOZIO
 Institut Mines-Telecom, Telecom ParisTech, CNRS LTCI, Paris, France

first.last@telecom-paristech.fr
 http://dbweb.enst.fr/

Context and Problem

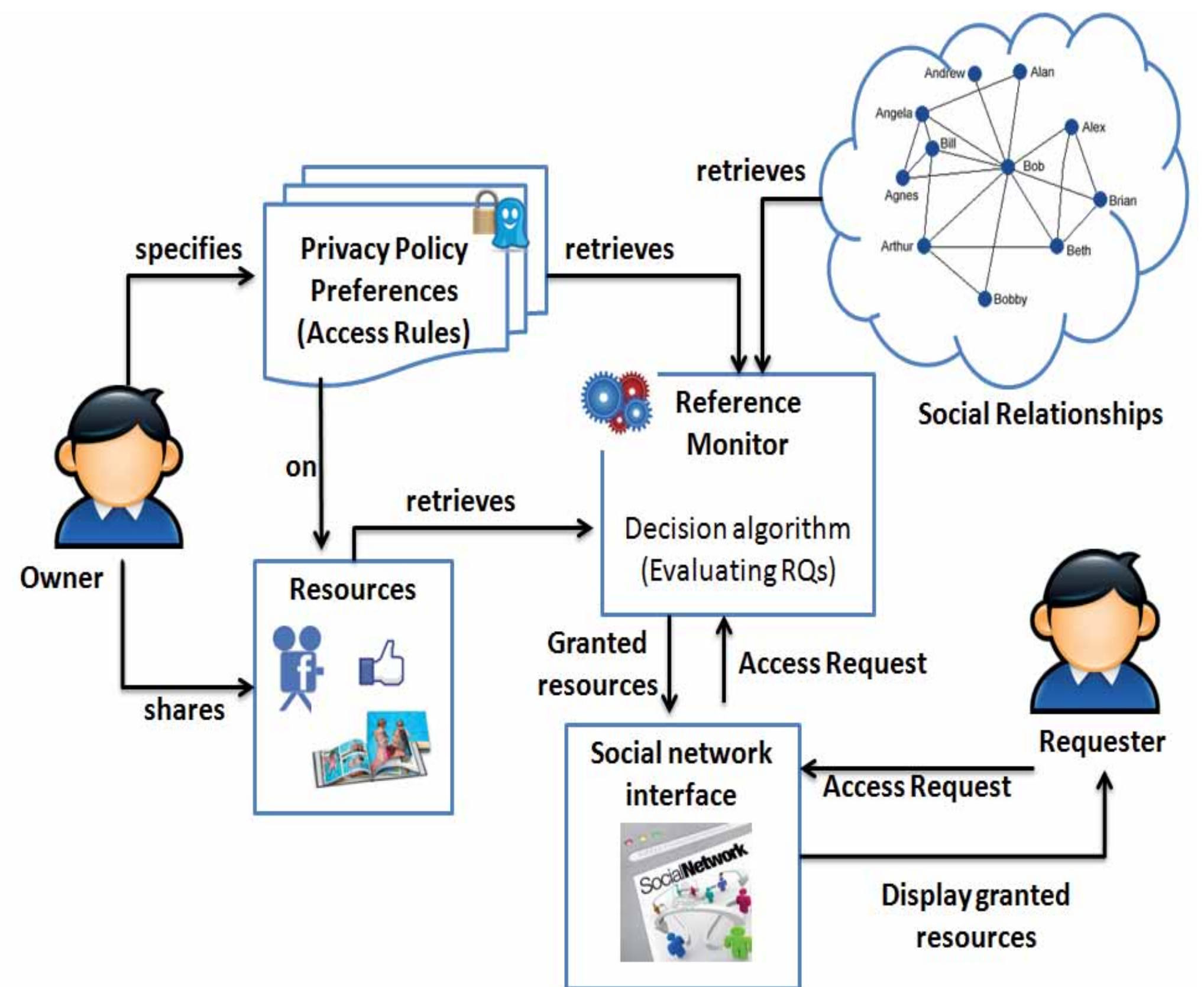
- Social network users :
 - have different kinds of relationships (friend, colleague, etc.)
 - share content (personal data, photos, videos, etc.)
- Social networks usually grow quickly in terms of number of users, relationships established and pieces of shared information.



Problem: Social network users have difficulty with specifying which information should be shared with whom.

Solution : Enable users to specify their privacy policy preferences in a more **flexible** and **efficient** way than existing privacy management systems.

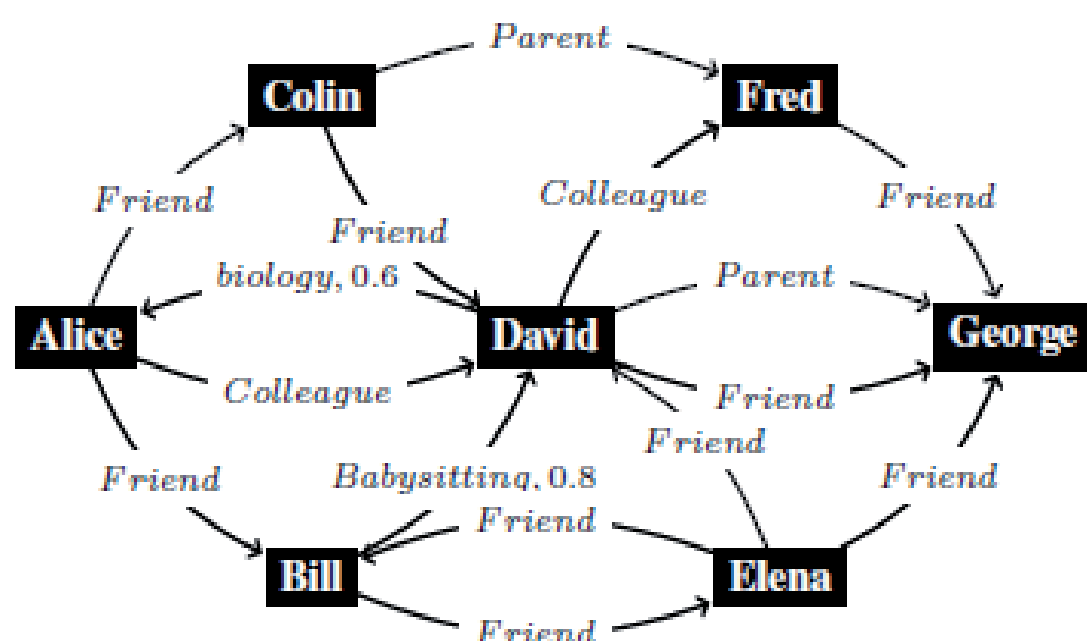
System Overview



Access Control Model

- **Social network model :**
 - Directed, edge-labeled, and weighted graph

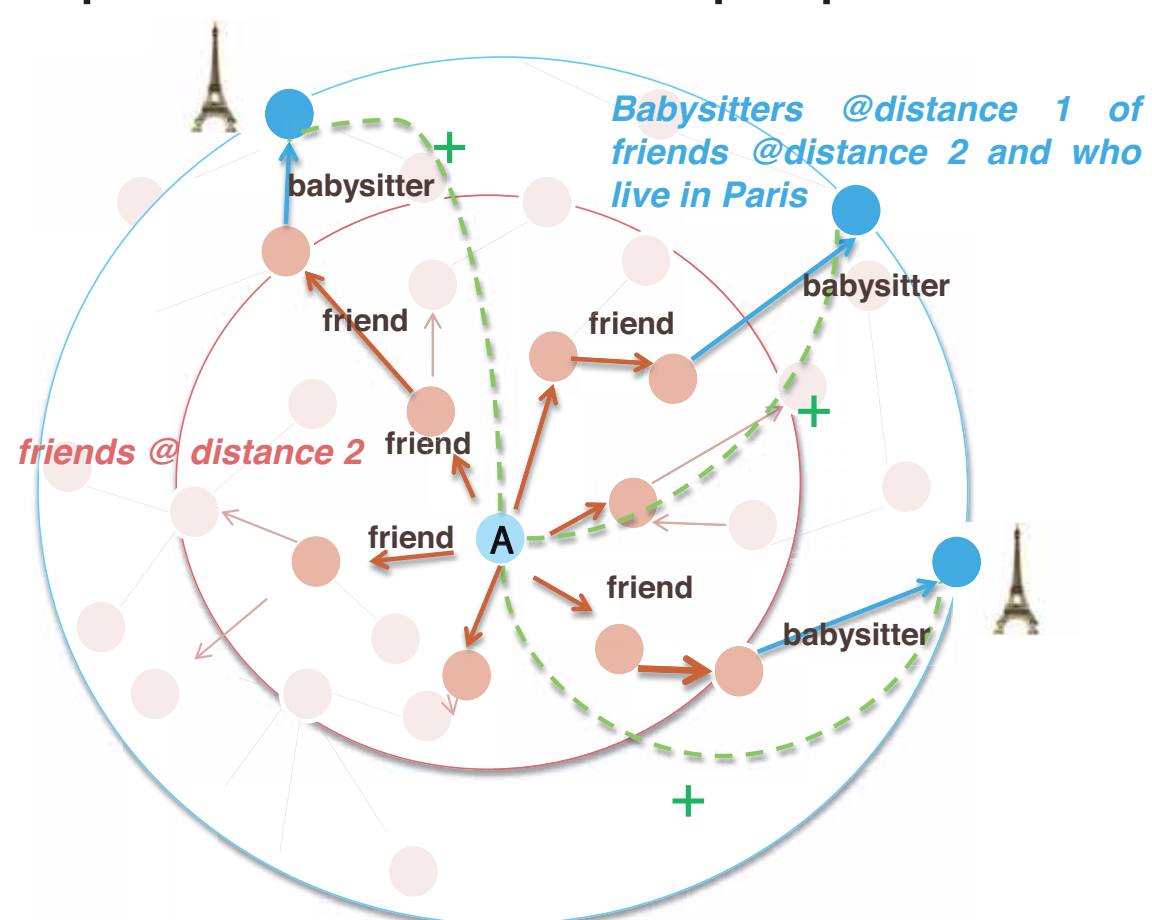
$$(V, E, \Sigma, \varphi)$$



A Social Network Subgraph

- **Access rules specification based on reachability constraints :**

- Semantics of the links, Links direction, Indirect relationships, Distance, User properties, Trust.



- **Access Rule (AR) :** Specification of the profiles of authorized users to access a given resource.

$$AR = (o, p, t_{min})$$

$$p = \{s_i\}$$

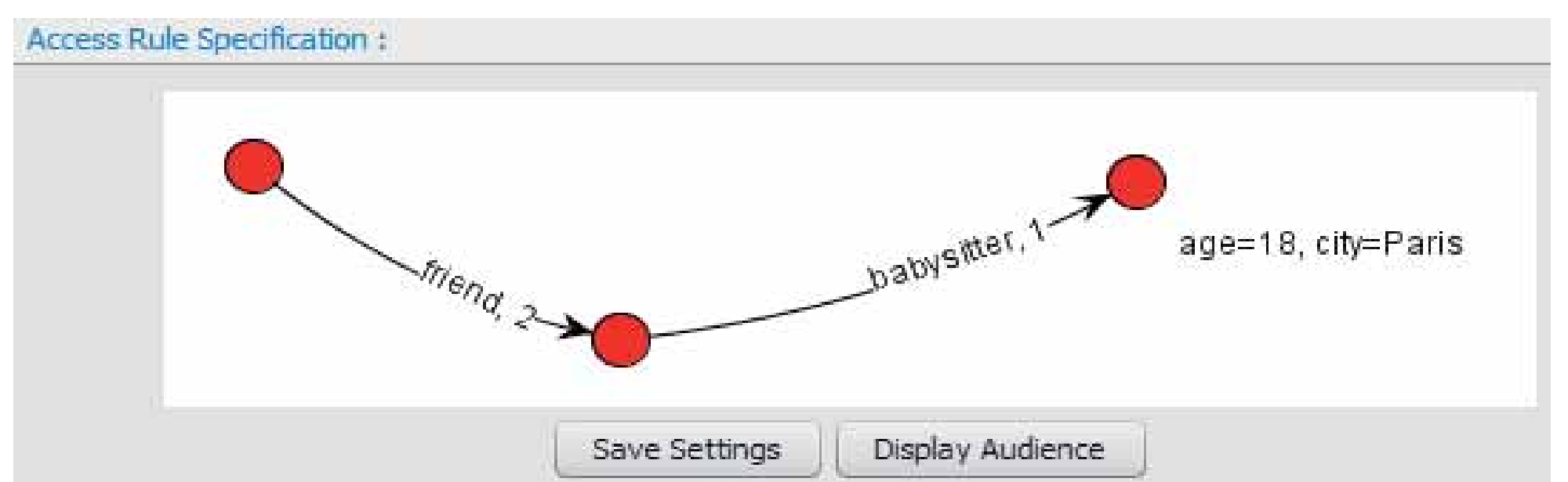
$$s_i = (r, dir, l, C)$$

- **Example :**
 - $p = \text{Friend} + [1, 2][\text{city} = \text{Paris}]/\text{BabySitter} + [1][\text{city} = \text{Paris}]$
 - $t_{min} = 0.8$

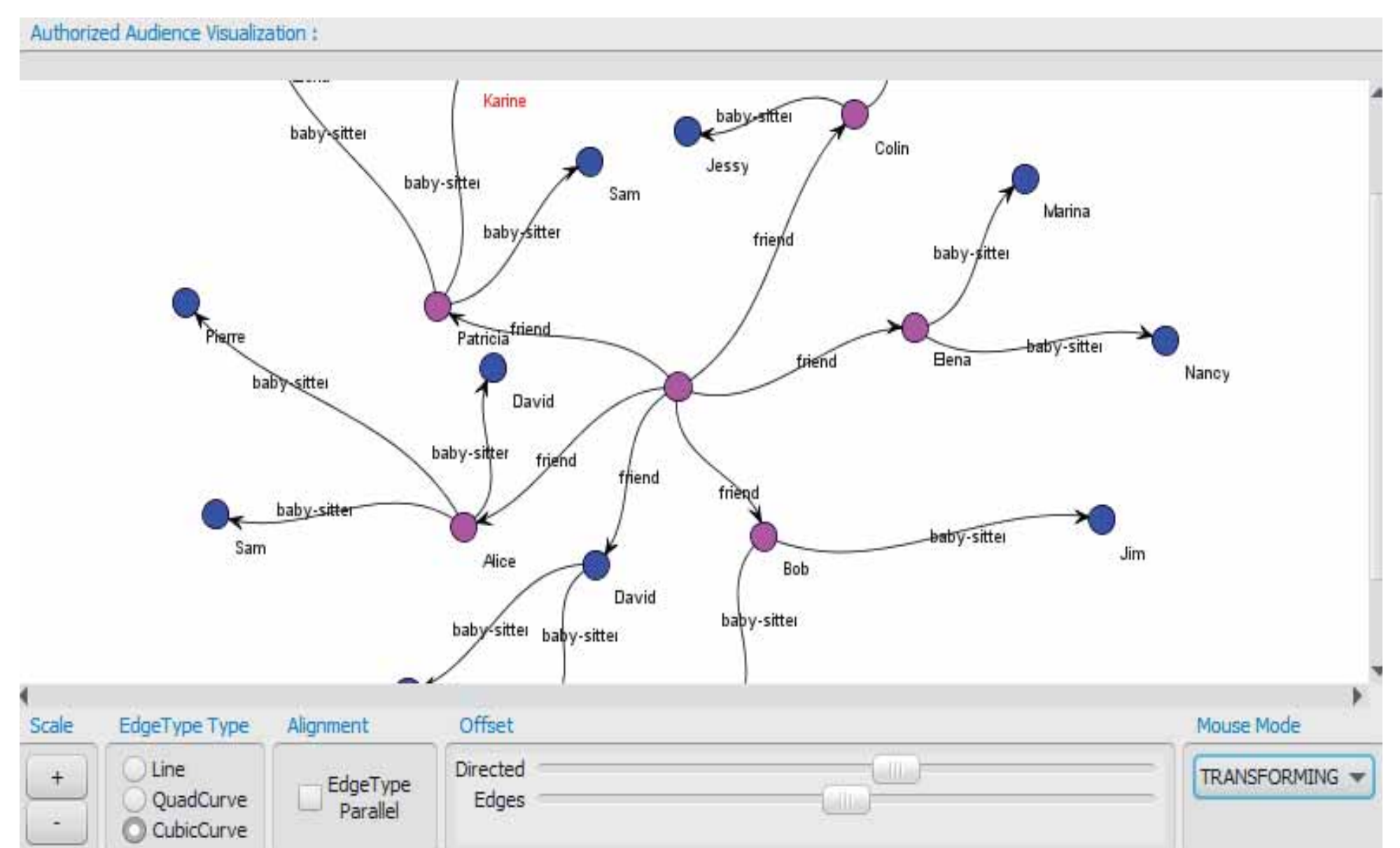
Demonstration

- **Dataset :** LiveJournal social network, ~5 Million users, ~80 Million edges

- **Access Rule Specification :**



- **Authorized Audience Visualization :**



A Probabilistic XML Merging Tool



Talel Abdessalem
Télécom ParisTech
Paris, France

Mouhamadou Lamine BA
Université Cheikh Anta DIOP
Dakar, Senegal

Pierre Senellart
Télécom ParisTech
Paris, France



<http://dbweb.enst.fr/>

What this tool aim at...

- Representing the outcome of semi-structured documents integration as a probabilistic tree
- Evaluating the uncertainty (modeled as probability values) of the result of the merge
- Querying the probabilistic repository with a subset of the XPath query language

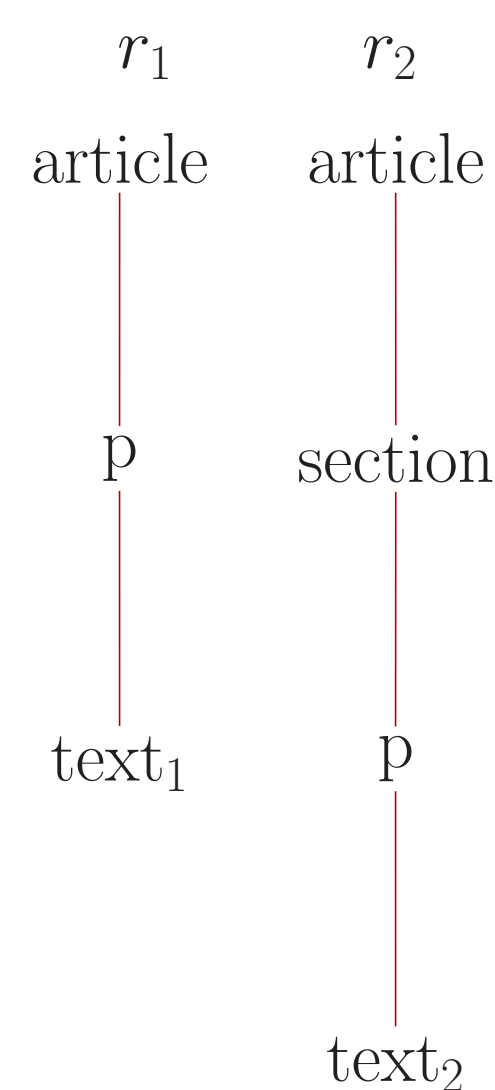
Application domain: [Wikipedia revisions](#)

The tool enables merging the revisions of a given Wikipedia page with:

- an efficient evaluation of the uncertainty of the obtained result
- an automatic management of conflicts.

Merging of Wikipedia revisions

- A two-way tree merging technique for P-Documents
- Two steps: [Matching of Revisions](#) and [Merging Matches](#)



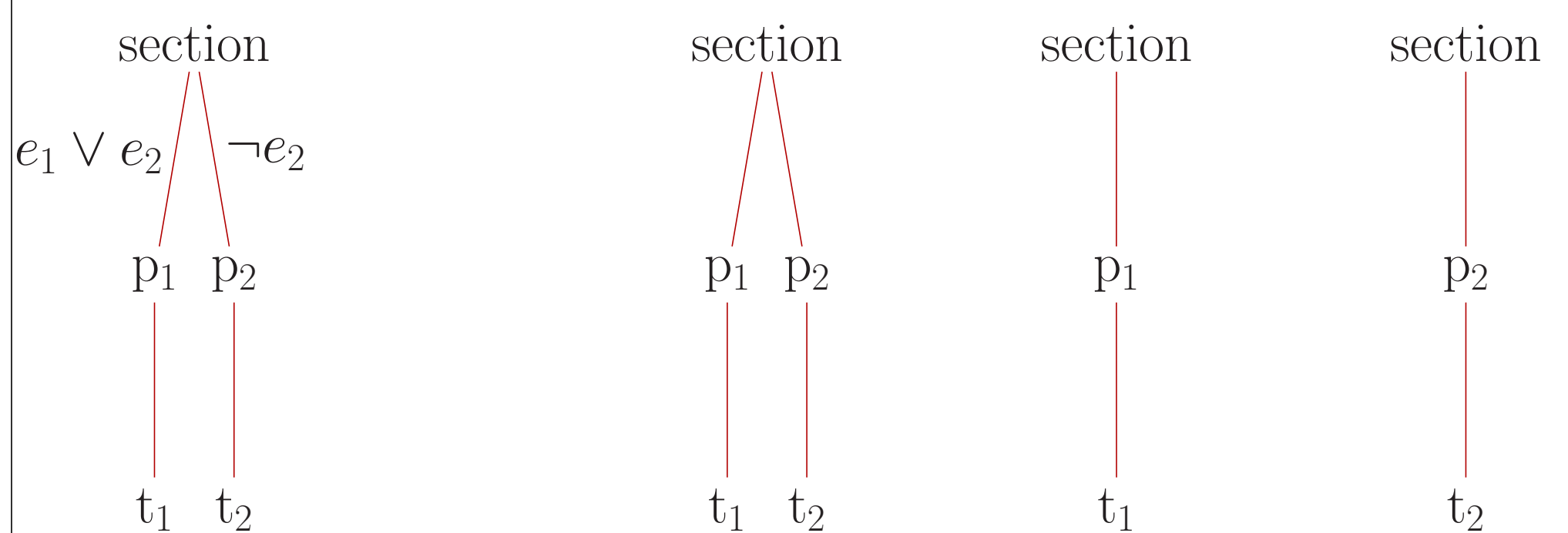
1. Matching of Revisions

Input: two revisions r_{k-1} and r_k and their associated event formula.

Output:

- Deleted nodes x : $x \in r_{k-1}$ and x has no match in r_k .
- Added nodes x : $x \in r_k$ and x has no match in r_{k-1} .
- Matched couples (x, y) : $x \in r_{k-1}$ and $y \in r_k$ match.

Probabilistic Documents



$\Pr(e_1) = 0.7$

$\Pr(e_2) = 0.6$

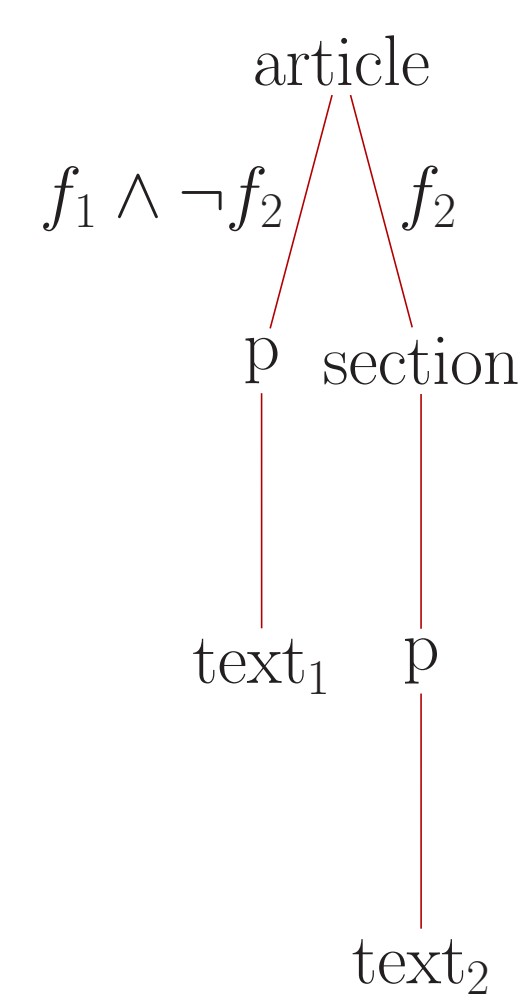
P-Document

0.28

0.6

0.12

Corresponding possible documents and their probabilities



The result of the merge process

2. Merging Matches

- Deleted nodes:

$$fie_{new}(x) = fie_{old}(x) \wedge (\neg f_k)$$

- Matched couple:

$$fie_{new}(x) = fie_{old}(x)$$

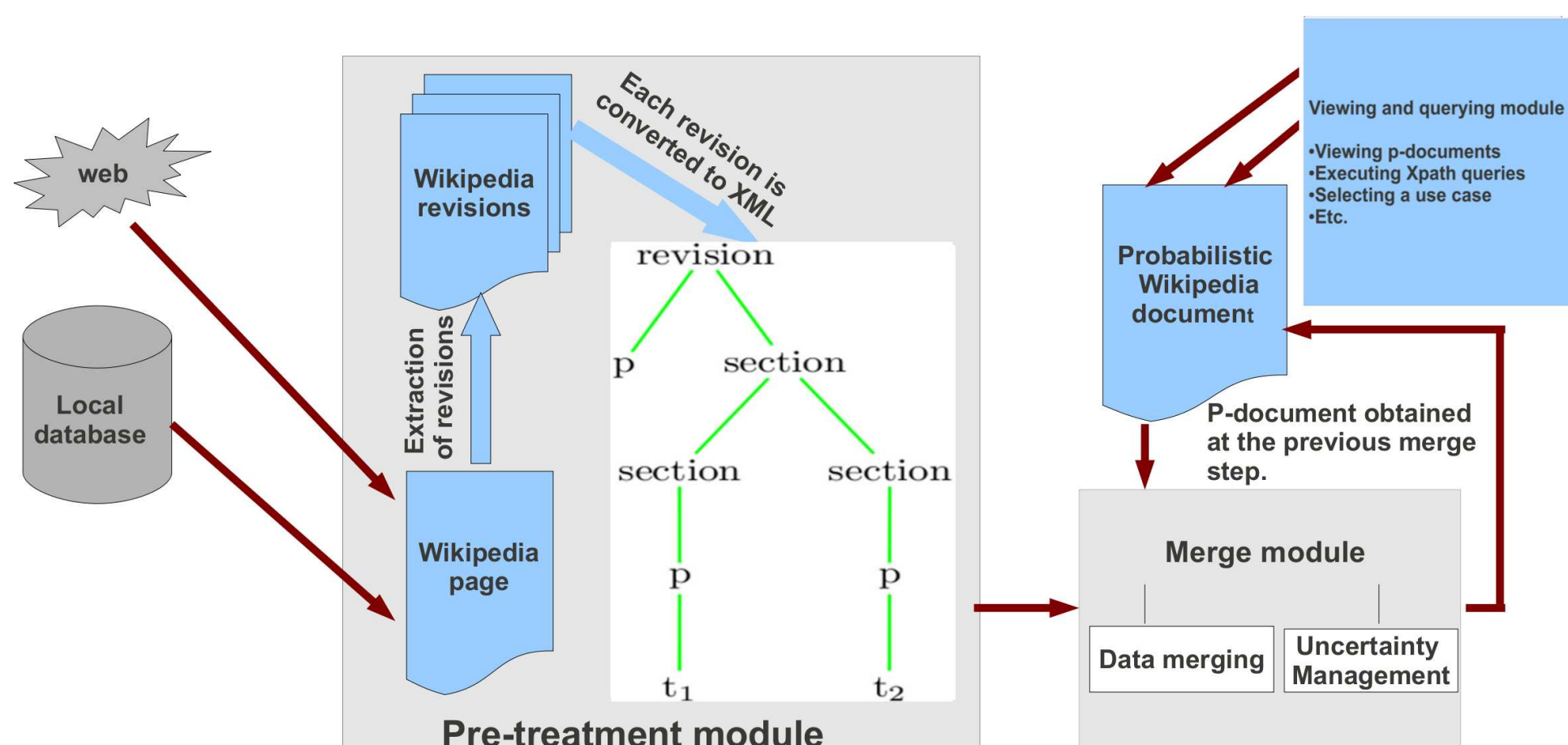
- For added nodes:

$$fie_{new}(x) = f_k$$

or

$$fie_{new}(x) = fie_{old}(x) \vee f_k$$

Architecture of the system



Description of the system

- [System for managing Wikipedia documents.](#)

Features

- A keyword-based search engine for Wikipedia pages
- Extracting the revisions of a given page
- Selecting the list of revisions to merge
- Building one's own Wikipedia article
- Displaying the result of the merge
- Demonstrating a certain number of use cases
- Using a subset of XPath query language

Parties prenantes



Auteurs

S. Cléménçon,
P. Bianchi,
G. Morral et
J. Jakubowicz

BACKGROUND

Motivation and applications

- Problem : investigate the **binary classification**
- Context : processing **BigData** for **statistical learning**
- Solution : implement in an **on-line** and **distributed** fashion

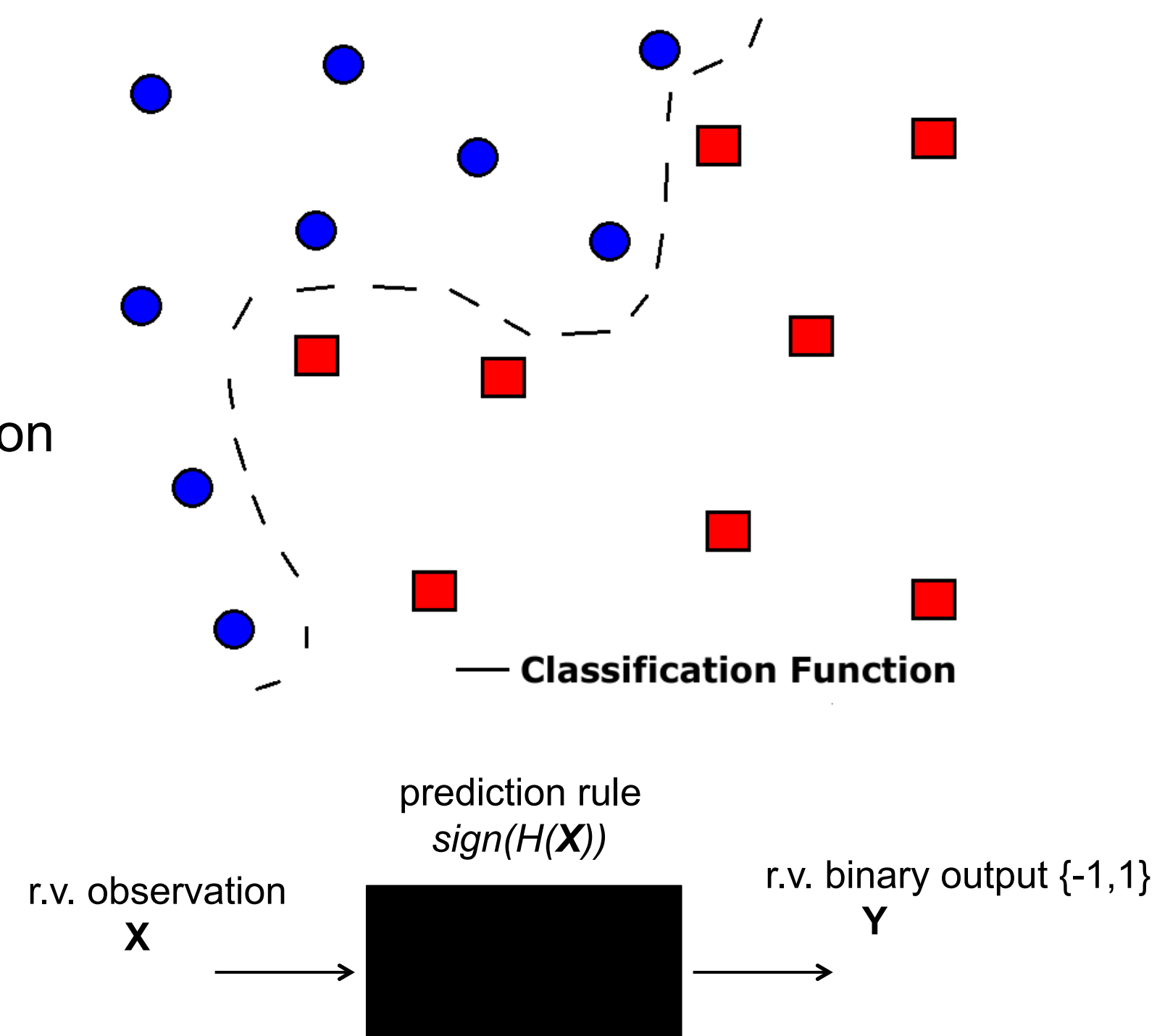
Learning problem

Given a training data set $(\mathbf{X}, \mathbf{Y}) = (X_i, Y_i)_{i \geq 1}$, find the best prediction rule $\text{sign}(H^*)$ such the classifier function $H(\mathbf{X}, \boldsymbol{\theta})$:

$$\boldsymbol{\theta}^* = \min_{\boldsymbol{\theta}} R_{\varphi}(H(\mathbf{X}, \boldsymbol{\theta})) \rightarrow \text{minimizes the Risk function } R_{\varphi}$$

Particular case :

- quadratic cost φ : $R_{\varphi}(H(\mathbf{X}, \boldsymbol{\theta})) = \frac{1}{2} E[(1 - YH(\mathbf{X}, \boldsymbol{\theta}))^2]$
- mixture of experts : $H(\mathbf{X}, \boldsymbol{\theta}) = \sum_j \theta_j h_j(\mathbf{X})$



PROPOSED DISTRIBUTED LEARNING

On-line Learning Gossip Algorithm (OLGA)

- A **distributed stochastic gradient descent** approach where the estimated parameter sequence $(\boldsymbol{\theta}_n)_{n \geq 1}$ is performed in 2 steps :

[Gossip step] At iteration n , each agent i transmits $X_{n,i}$ to all randomly selected neighbours j with *probability* p and obtain $h_j(X_{n,i}, \boldsymbol{\theta}_{n-1,j})$

[Local descent step] each agent i update its estimated parameter $\boldsymbol{\theta}_{n,i}$ as follows :

$$\boldsymbol{\theta}_{n,i} = \boldsymbol{\theta}_{n-1,i} + \gamma_n \nabla_i h_i(X_{n,i}, \boldsymbol{\theta}_{n-1,i}) (Y_{n,i} - Y'_{n,i}^{(V)})$$

where

- $(\gamma_n)_{n \geq 1}$ is a decreasing step size sequence
- the r.v. $Y'_{n,i}^{(V)}$ is an unbiased estimate of the global decision $H(X_{n,i}, \boldsymbol{\theta}_{n-1,j})$ given by $Y'_{n,i}^{(V)} = h_i(X_{n,i}, \boldsymbol{\theta}_{n-1,i}) + 1/p \sum_j \delta_{n,i}^j h_j(X_{n,i}, \boldsymbol{\theta}_{n-1,j})$ and $(\delta_{n,i}^j)$ are independent Bernoulli r.v 's $B(p)$

Asymptotic behaviour

Theoretical results (under suitable assumptions)

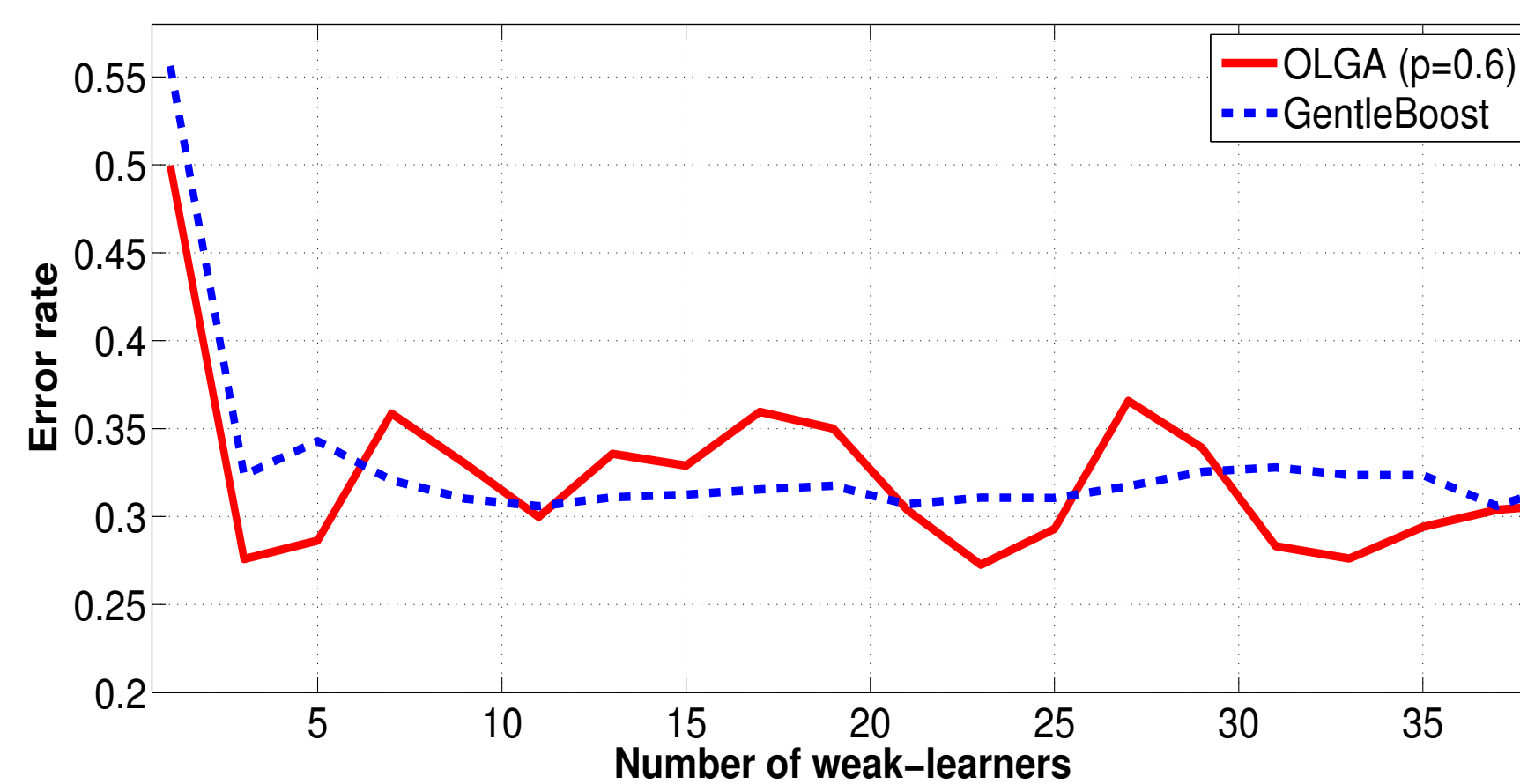
✓ Consistency : $(\boldsymbol{\theta}_n)_{n \geq 1}$ **convergence a.s.** to the set of stationary points of R_{φ}

✓ Conditional Central Limit Theorem : qualify the error variance excess

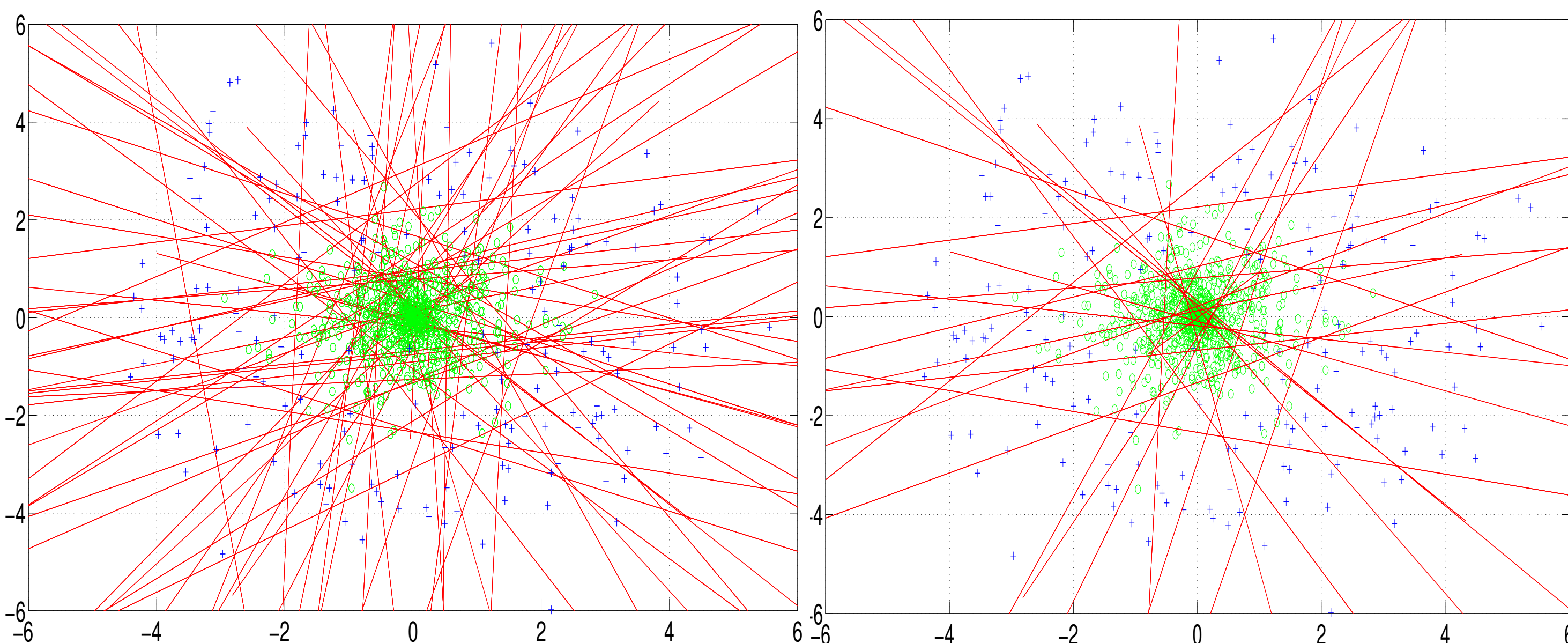
$$\gamma_n^{-1/2} (\boldsymbol{\theta}_n - \boldsymbol{\theta}^*) \rightarrow N(0, \Sigma(\boldsymbol{\Gamma}^*))$$

$\boldsymbol{\Gamma}^*$: error in a centralized case + **error excess induced by sparsification**

✓ The average network throughput is **reduced** by a factor **(1-p)**



Performance comparison between the centralized GentleBoost and OLGA for a benchmark dataset



Result classification with OLGA (-) on a simulated binary dataset (+ and o) using weak classifiers

Left : OLGA
Right : OLGA with **agent selection**

at each iteration n , each agent i of V declares idle or **active** under a suitable criterion \rightarrow time-varying agents set V_n

- ✓ Reduce redundancy classifiers and keep the relevants



Demonstrating Intelligent Crawling and Archiving of Web Applications



arcomem

Muhammad Faheem
Institut Mines-Télécom
Télécom ParisTech; CNRS LTCI
Paris, France

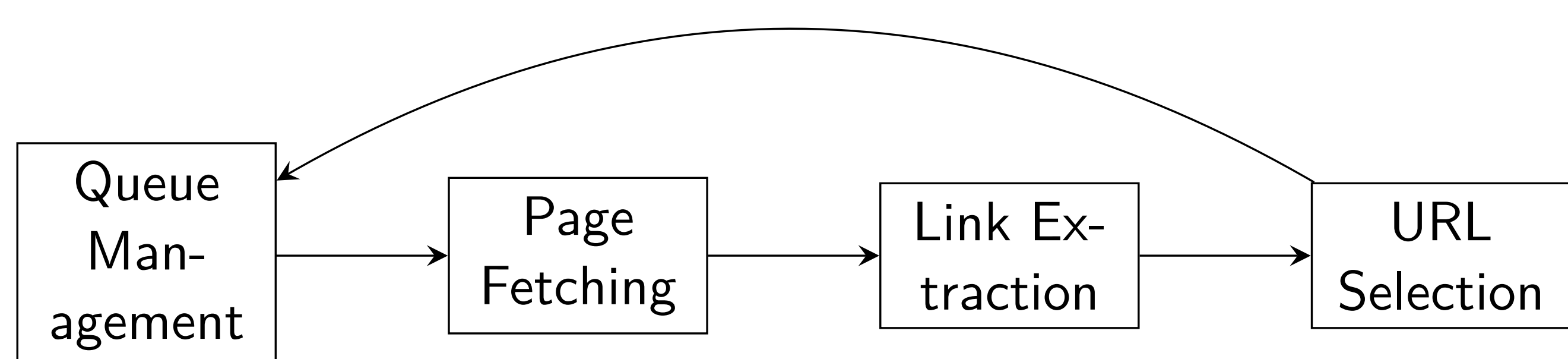
muhammad.faheem@telecom-paristech.fr

Pierre Senellart
Télécom ParisTech
& The University of Hong Kong
Hong Kong

pierre.senellart@telecom-paristech.fr

Traditional crawler

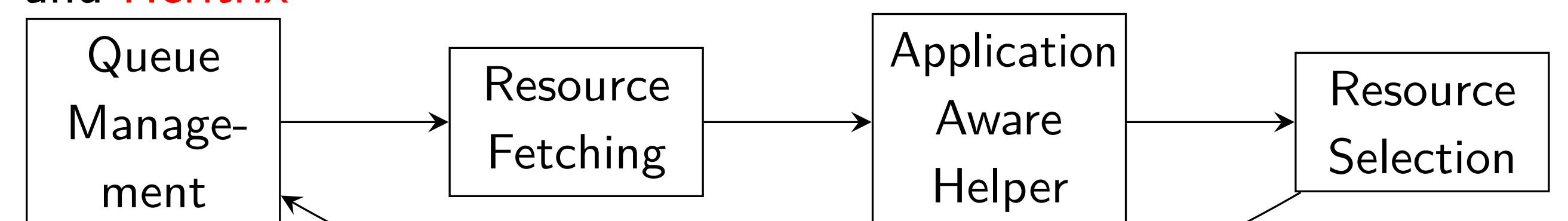
Traditional crawling: independent of the nature of the sites and their content management system



⇒ Many HTTP requests, no guarantee of content quality

Application-aware helper

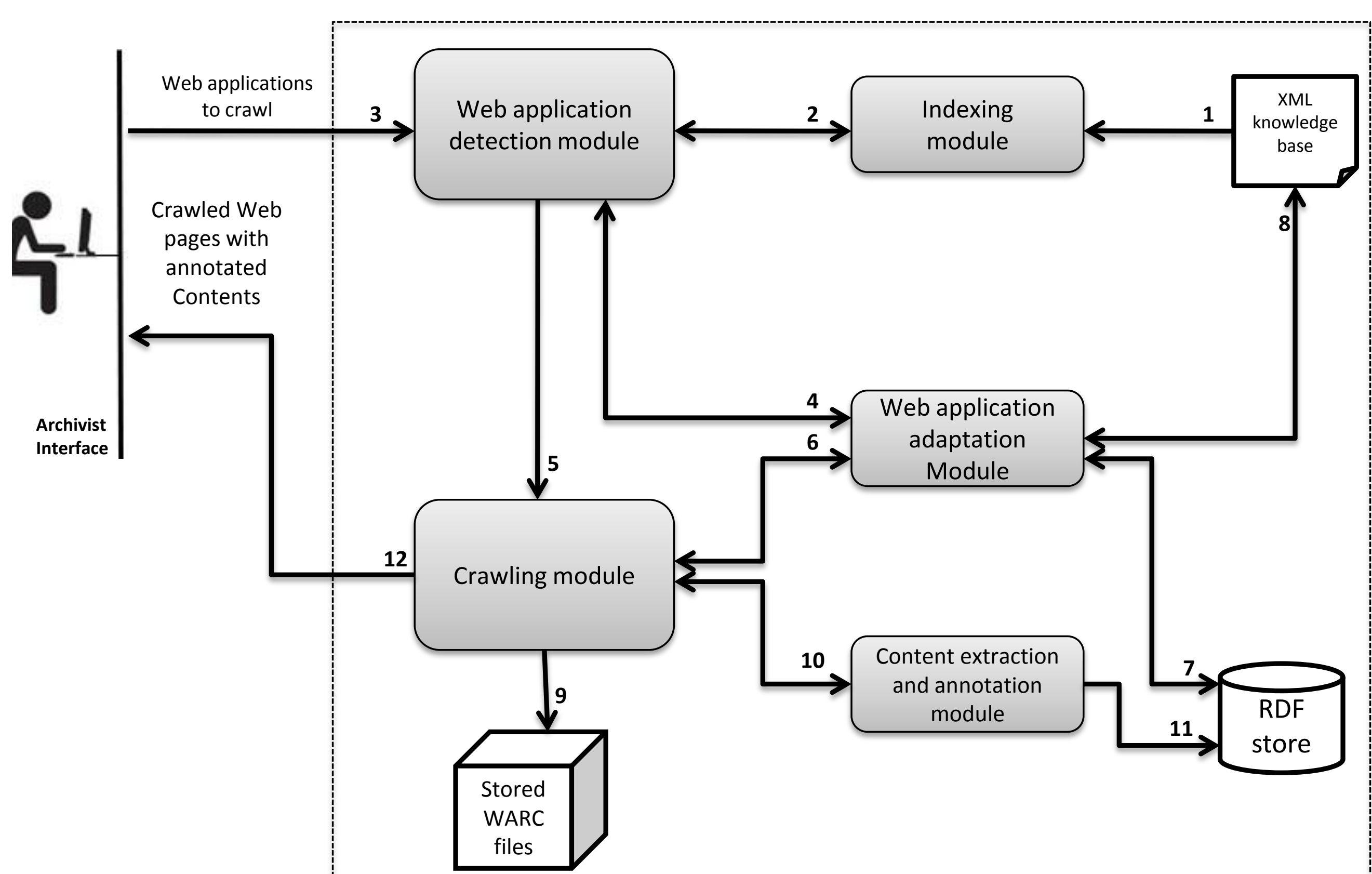
- Different crawling techniques for different Web sites
- Detect the type of Web application, kind of Web pages inside this Web application, and decide crawling actions accordingly
- Directly targets useful content-rich areas, avoids archive redundancy, and enriches the archive with semantic description of the content
- Implemented in 2 Web crawlers: Internet Memory Foundation crawler and Heritrix



Goal: Smart archiving of the Social Web:

1. Performing intelligent Crawling
2. Archiving Web objects

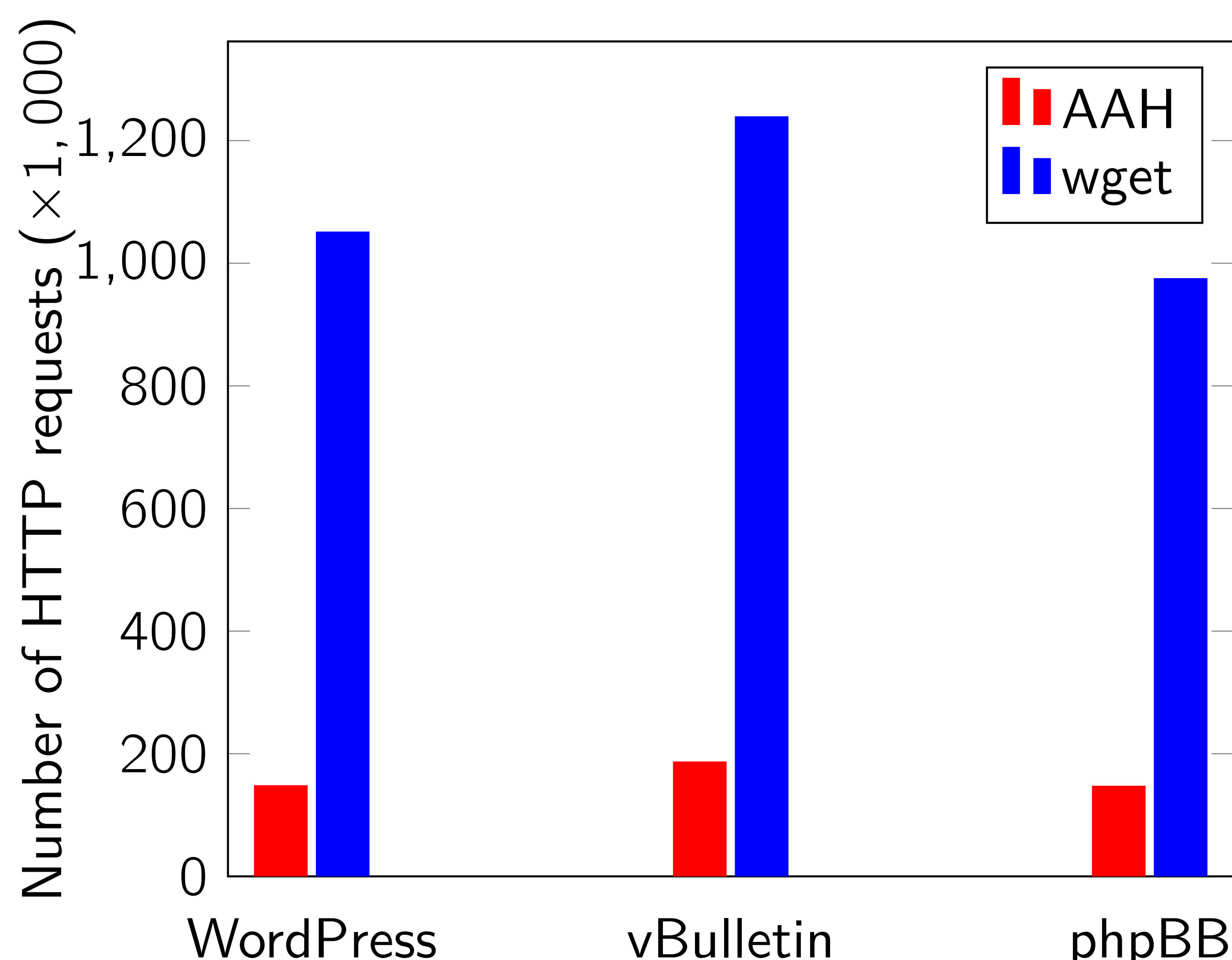
Architecture



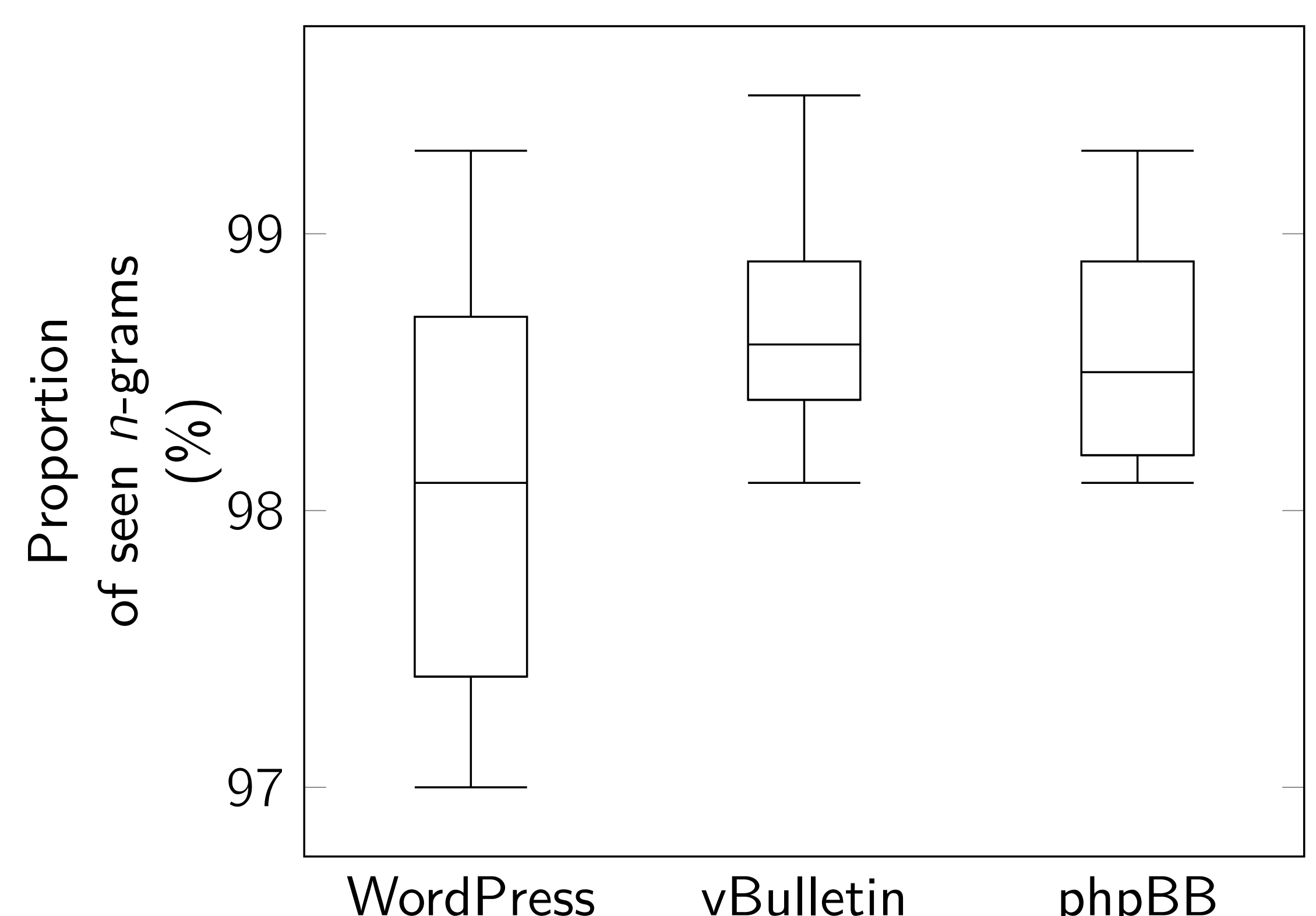
Methodology

- Knowledge base of known Web application types, algorithms for flexible and adaptive matching of Web applications to these types
Declarative, XML-based format
Integrated with YFilter for efficient indexing of KB.
- Type detected using URL patterns, HTTP metadata, textual content, XPath patterns, etc. E.g., `vBulletin Web forum: contains(//script/@src, 'vbulletin_global.js')`
- Different crawling actions for different kinds of Web pages under a specific Web application
- Crawling action: not just a list of URLs; can be any action that uses REST API, complicated interaction with AJAX-based application, and extracts semantic Web objects

Crawl efficiency



Crawl effectiveness



CrowdMiner: Mining association rules from the crowd

Introduction

- **Crowd data sourcing** collects data from the crowd, often by asking questions
- We want to learn about new domains from the crowd
 - E.g., health-related habits in some population
- Data is not recorded anywhere
- The contents of the domain are unknown
 - Discover what is **interesting** about this domain

What should we ask the crowd?

Data mining for the crowd?

- The discovery of data patterns in databases is done by **data mining**.
 - Not suitable for our case
 - People do not remember enough details!
- For example, it is unrealistic to expect people to remember every activity they did in the past, everything they have eaten, etc.
- They are far more likely to remember **personally prominent patterns**

"I drink red wine about once a week"

The model

We learn *association rules* of the form $a, b \rightarrow c, d$

- E.g., "heartburn" \rightarrow "baking soda", "lemon"

The answers contain

- **Rule support** – frequency of a, b, c, d
- **Rule confidence** – frequency of c, d given a, b
- **Items** (for an open question)
- **Significant rules** – average user support and confidence exceed fixed thresholds
- Users treated as random samples

Our approach

- Use **personal summaries** to learn about **general trends**
- Treat individual answers as samples
- Combine two types of questions

- **Open questions**

"Complete: When I feel tired, I usually go for a walk"

- **Closed questions**

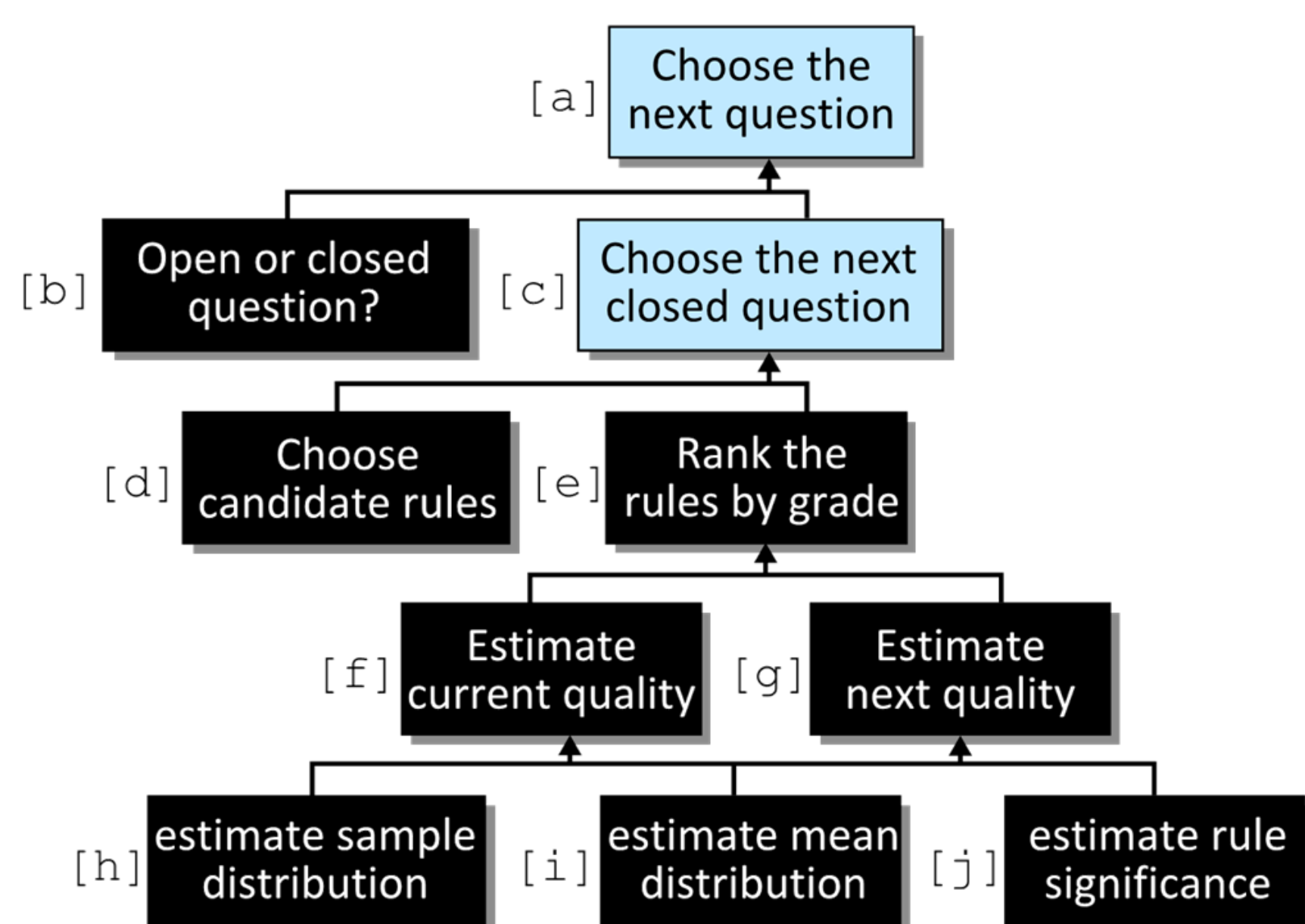
"When you have a heartburn, do you take baking soda and lemon?"

- Easier for users to answer
- Help digging deeper into their memories

We develop a system prototype *CrowdMiner* that interactively decides what to ask in order to discover significant data patterns

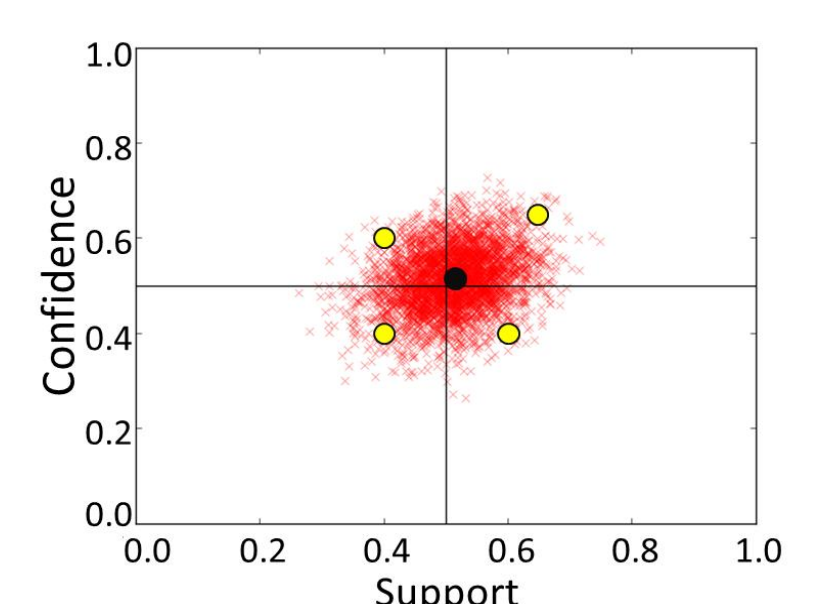
Choosing the Questions

A hierarchy of components that allow estimating the effect of the next question and choosing accordingly



Error Estimations

- Not all the users can be asked about every rule
- We want to estimate the probability of making an error – given the current knowledge
 - We learn a distribution of the answer support and confidence
 - **Significance estimation** – by the position of >0.5 of the distribution mass
 - **Error probability** – for the true mean to be on the other side of the thresholds
- The next question is the one expected to minimize the overall error

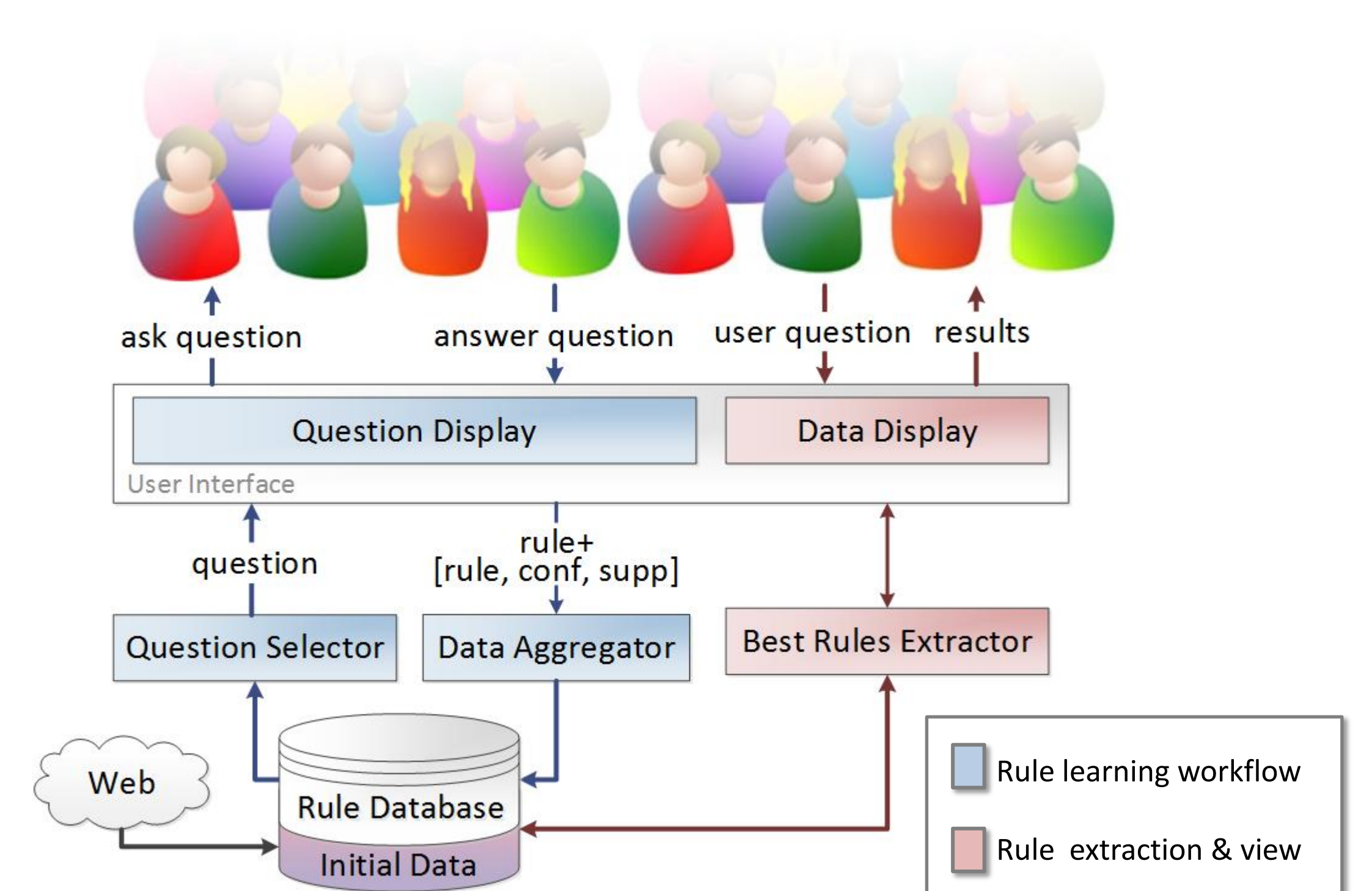


Well-Being Portal

- Learn about the **health habits** of others – by browsing the portal
 - Sports activities, eating habits, natural treatments
 - ...
- Portal users are occasionally prompted with **questions**
 - About their personal habits
 - Computed by our algorithm
- User **answers** are processed to deduce rules (associations) between well-being concepts in the portal
- The portal allows browsing the learned rules



System Architecture



Institutions



Authors

Julio Cesar Louzada Pinto
Tijani Chahed
Jérémie Jakubowicz

Partners



Objectives

- Develop a stochastic opinion dynamics model with multiple contents and study the asymptotic behavior in simple cases.
- Develop a community detection (graph clustering) algorithm.

Model

- N agents communicate about K contents via a graph $G = (V, E)$ with inward adjacency matrix A .
- Agent i has score $X_t^{i,k}$ for content k at time t .
- Preferences $P_t^{i,k}$ are normalized scores, i.e., $P_t^{i,k} \propto X_t^{i,k}$.
- Agents update their scores linearly as

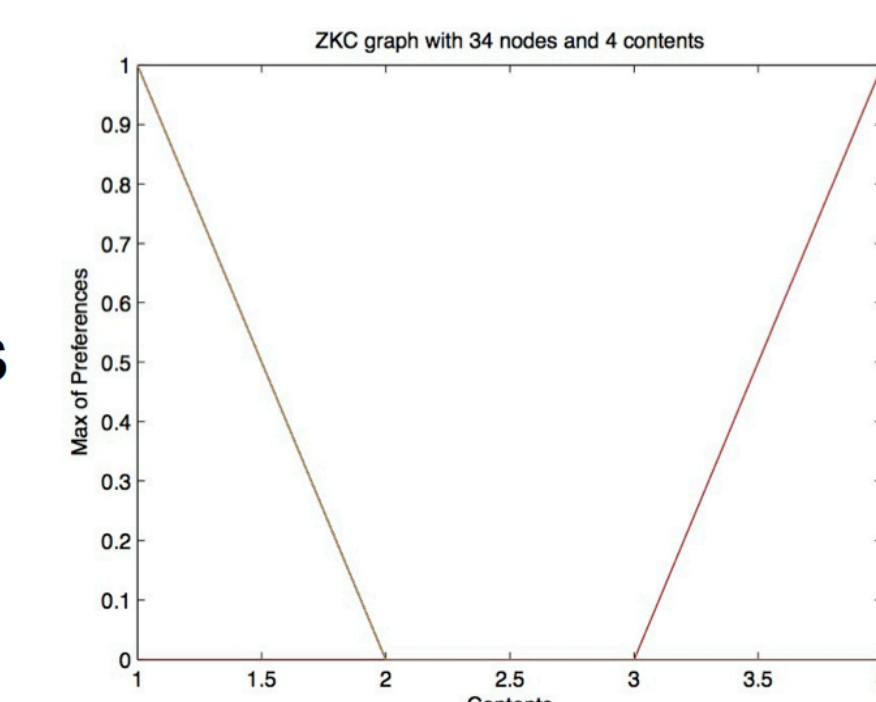
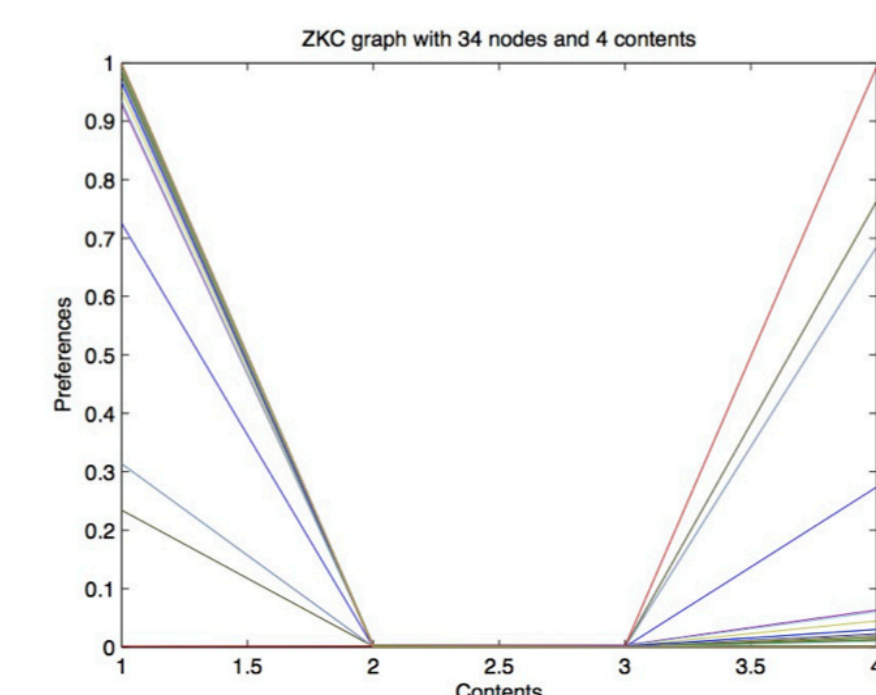
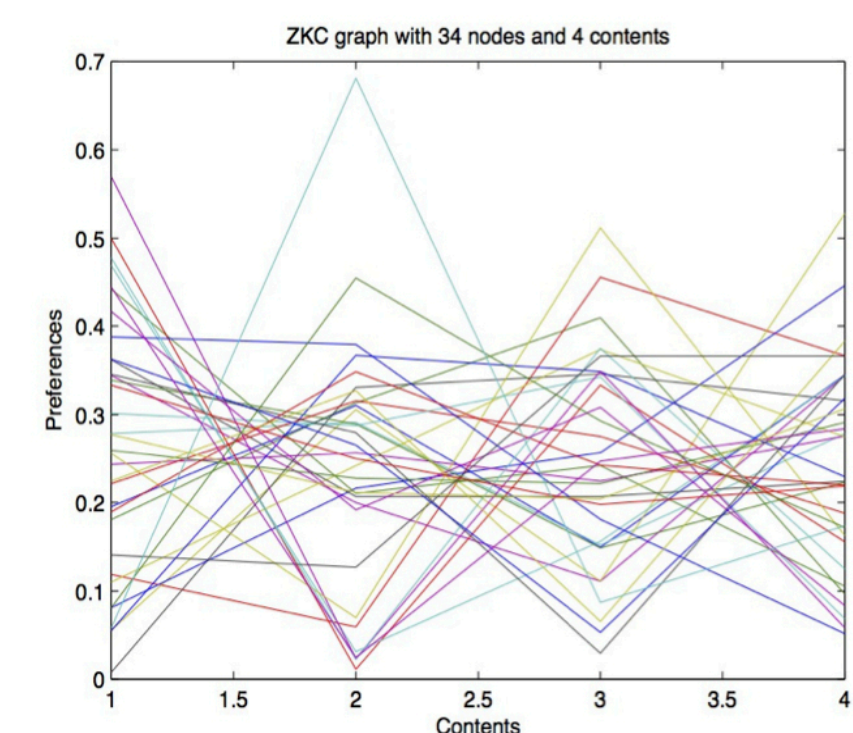
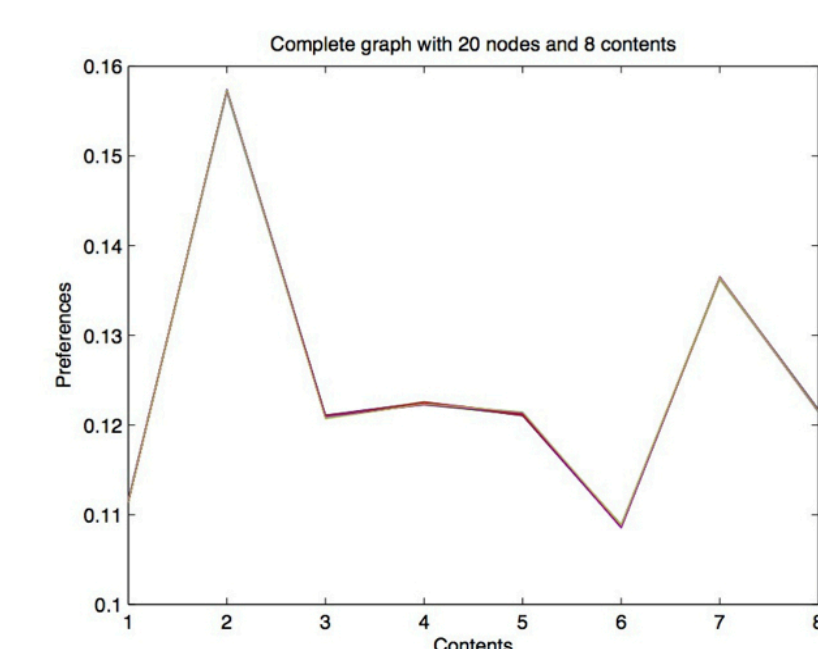
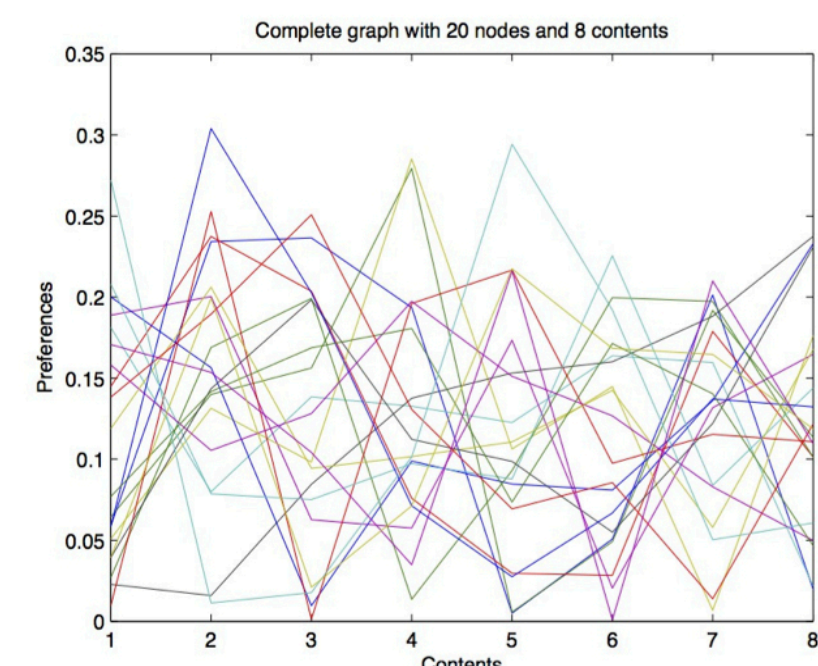
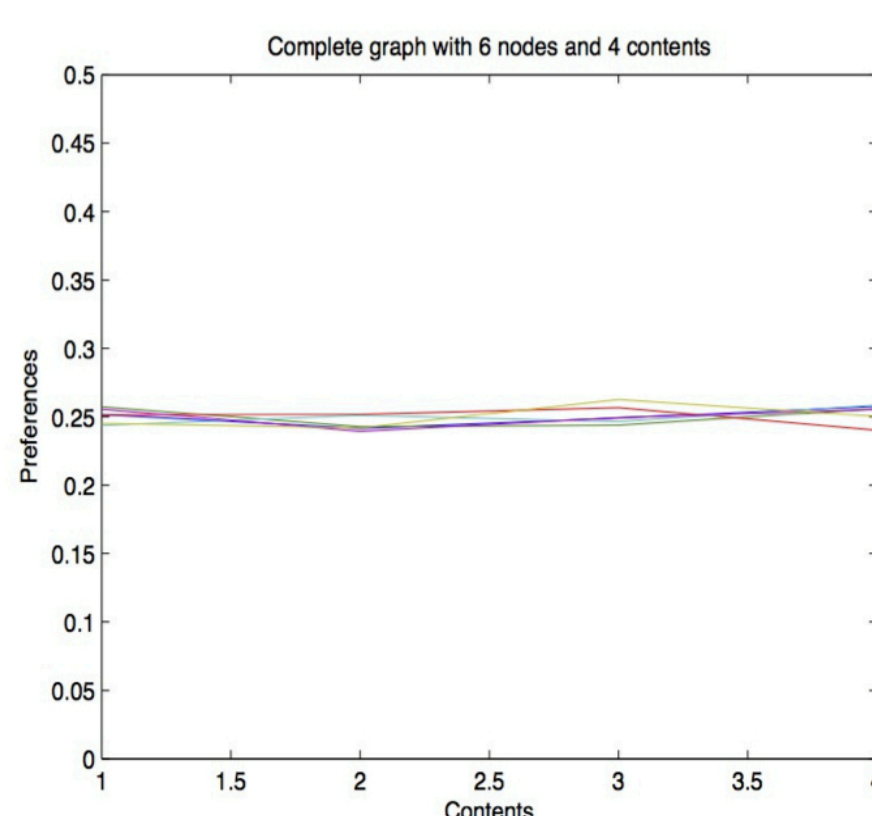
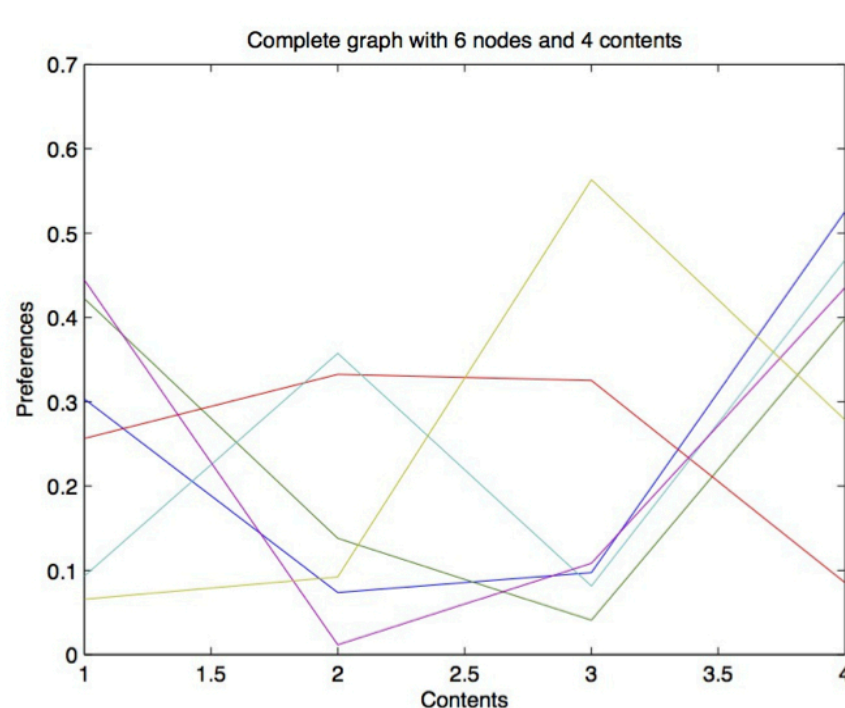
$$X_{t+1}^{i,k} = X_t^{i,k} + \sum_{j \rightarrow i} 1_{u_{t+1}^j = k}$$

with $P(u_{t+1}^j = k | F_t) = (f(P_t))_{jk}$ the probability of agent j broadcasting content k .

- Function $f: (\Delta_K)^N \rightarrow (\Delta_K)^N$ models the way agents choose the contents to broadcast, where $(\Delta_K)^N$ is the set of the $N \times K$ stochastic matrices.

Types of function f

- $f(x) = x$ – identity function.
- $(f(x))_{ik} = \frac{e^{\beta x_{ik}}}{\sum_{i'} e^{\beta x_{i'k}}}$ – soft-max function.
- $\beta \ll 1 \rightarrow (f(x))_{ik} \sim \frac{1}{K}$
- $\beta \gg 1 \rightarrow (f(x))_{ik} \sim 1_{k = \text{argmax}_j x_{ij}}$



Results

- For $f(x) = x$: there exists a random variable $P_\infty \in (\Delta_K)^N$ such that $P_t \rightarrow P_\infty$ almost surely.

If G is strongly connected, then there exists a random variable

$$\pi \in \Delta_K \text{ such that } P_t \rightarrow 1\pi^T \text{ almost surely.}$$

- For soft-max exponential with $\beta \ll 1$: If $\inf_i \sum_j A_{ij} > 0$, then there exists a $\beta_{min} > 0$ such that for all $\beta \in [0, \beta_{min}]$ we have that

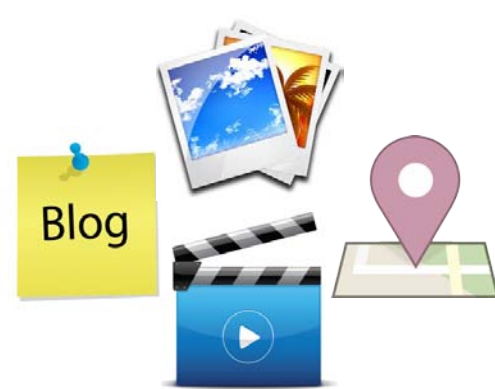
$$P_t \rightarrow \frac{11^T}{K} \text{ almost surely.}$$

- If $\beta \gg 1$, then we expect the graph to be clustered in communities which broadcast the same content; the one with the maximum preference. This creates an algorithm for community detection.
- Publication: A stochastic opinion dynamics model with multiple contents, CDC, Firenze, December 2013.

Authors : Daqing Zhang, Dingqi Yang, Zhu Wang, Zhiyong Yu

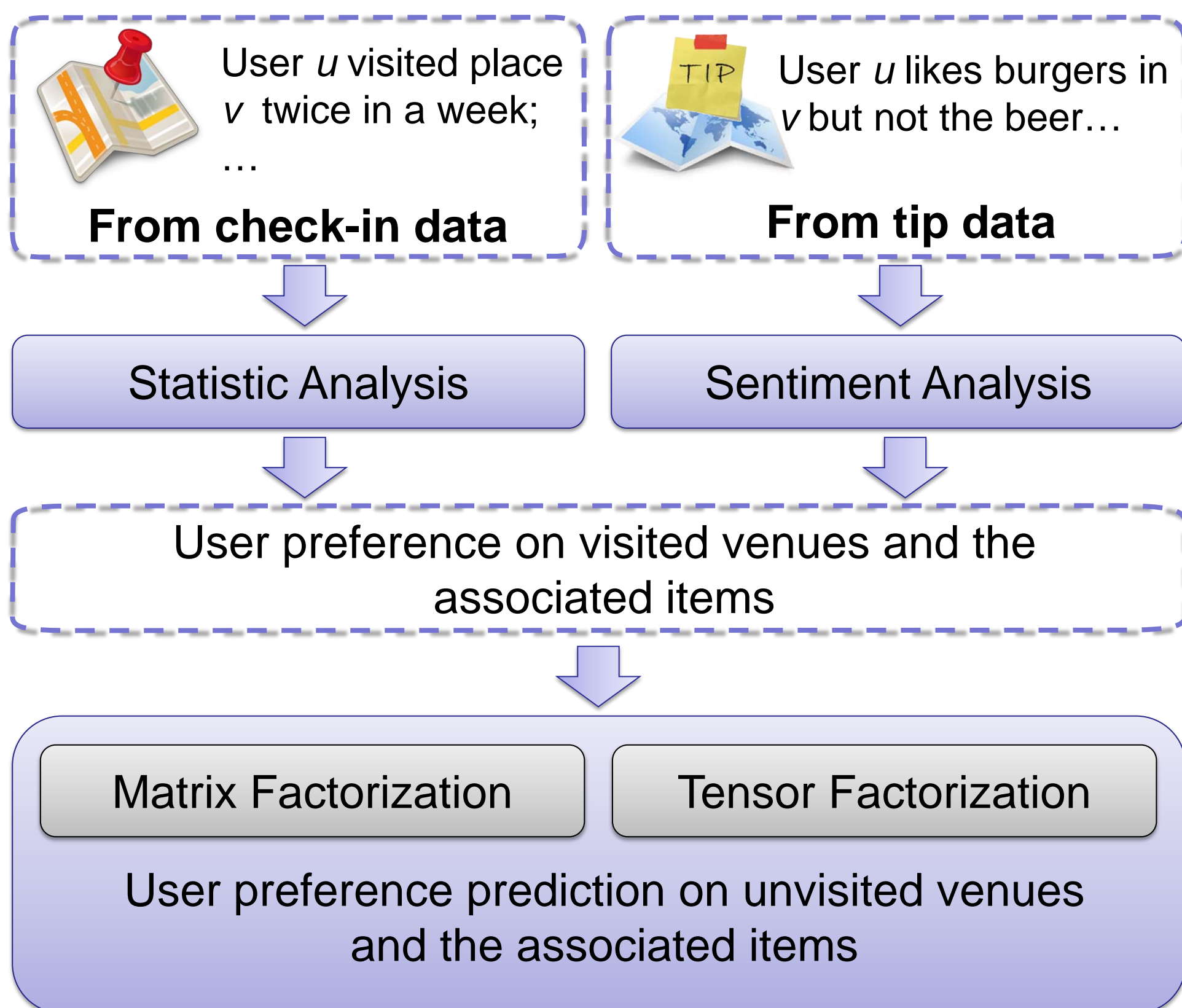
Location Based Social Network Data Analytics

- With the increasing popularity of location based social networks, users generated significant volume of heterogeneous social media, e.g.,
 - Texts
 - Photos
 - Videos
 - Presences
 - ...
- These digital footprints massively contain users' fine-grained preference
- Understanding this user preference can enable ubiquitous, personalized location based services.



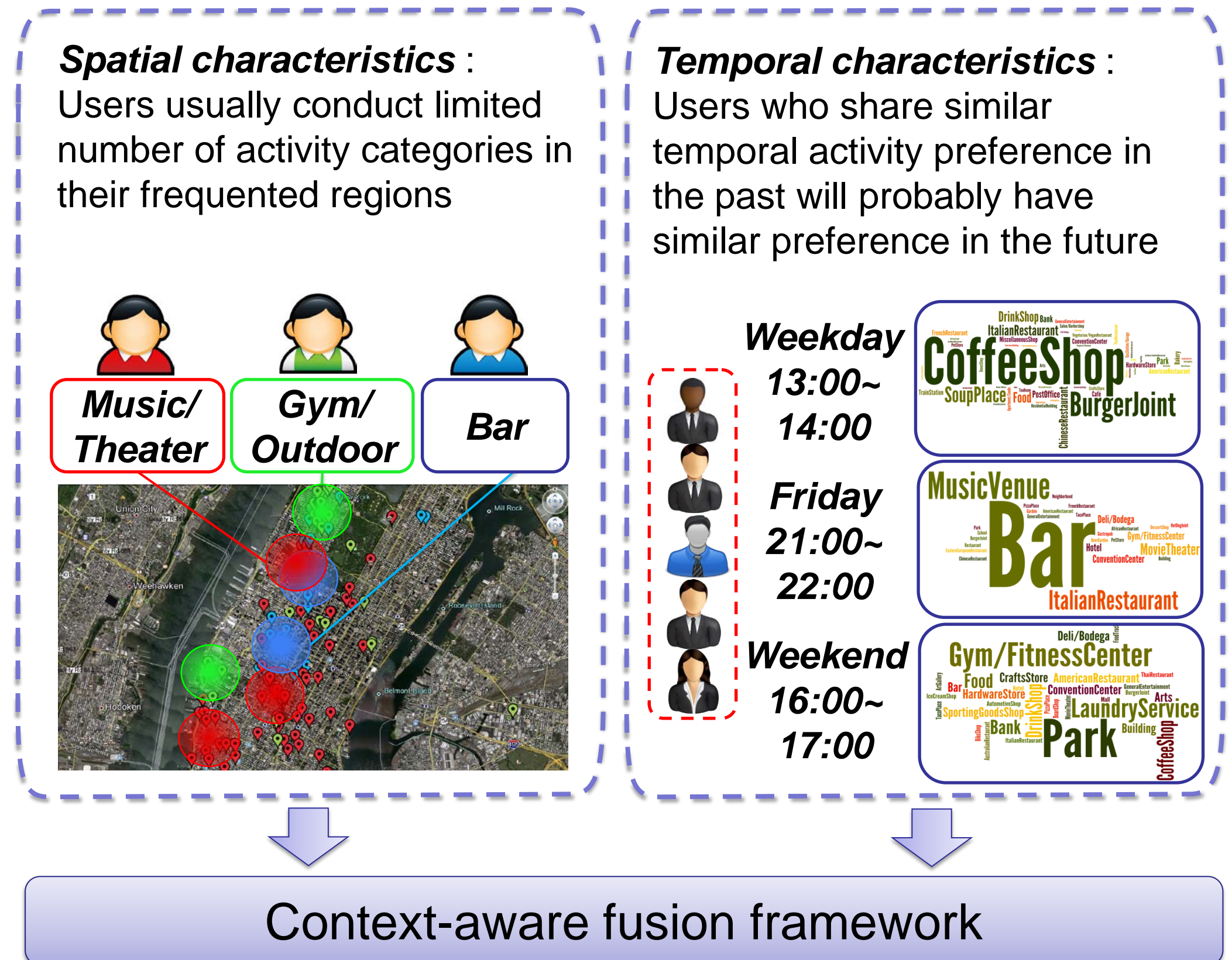
Preference Awareness

- Extracting fine-grained user preference on venues from heterogeneous data.
- Predicting user preference on unvisited venues.

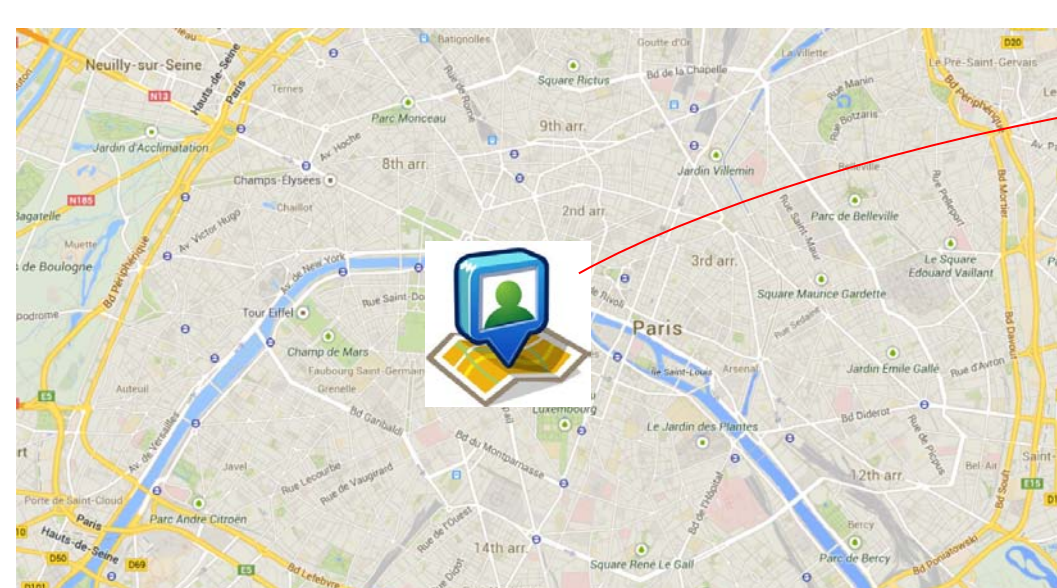


Context Awareness

- Studying and modeling spatial temporal characteristics of user activity.
- Inferring a user's interest according to his current context.

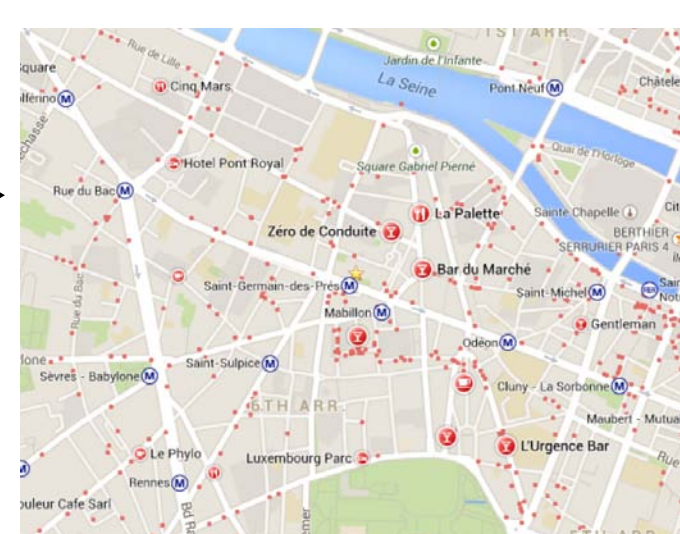


An Example of Personalized Location Based Services

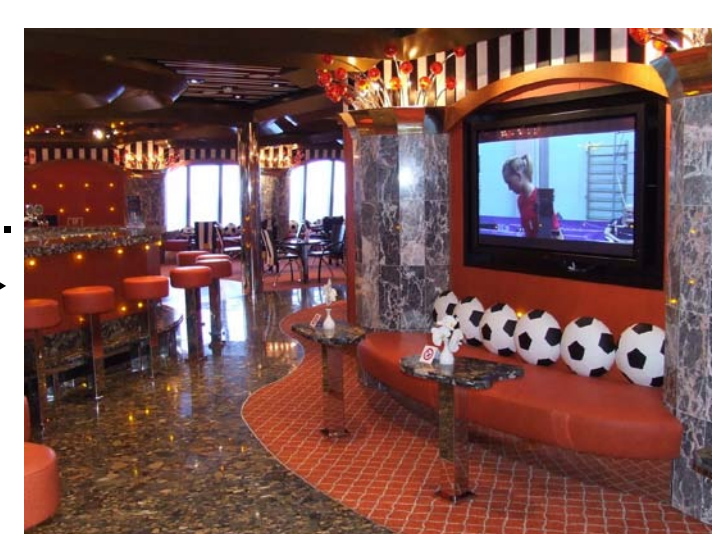


Time: 20:15 Friday
GPS: 48.8525,2.3344

He is probably interested in going to a bar.
Context-awareness



He will like the x bar and the large-screen TV and sportive environment there.
Preference-awareness



Authors

Djallel Bouneffouf,
Amel Bouzeghoub and
Alda Gançarski

Context



Objectives and Issues

Motivations

- Provide personalized and context-aware recommendations in mobile environments
- Consider content dynamicity and user's situations risk level

Key Challenges

- Infer higher-level goals from low-level observed operations
- Handle cold start and sparseness effect
 - Requires a large amount of information in order to make accurate recommendations
- Exploration vs. exploitation tradeoff
 - How to sacrifice a short term small reward to privilege larger rewards in the long term?
 - How to associate the situations risk level to the exploration/exploitation tradeoff?

Key Words & Key Technologies

- Context-Aware Recommender Systems (CARS)** combine characteristics from context-aware systems and recommender systems in order to provide personalized recommendations to users in ubiquitous environments.
- Machine learning** algorithms can be used to learn models and predict documents
 - Reinforcement learning** is learning what to do: how to map situations to actions
- In probability theory, the **multi-armed bandit problem** models an agent that simultaneously attempts to acquire new knowledge and optimize her decisions based on existing knowledge
 - In each round, a learner takes an **action** (or **arm**) and in return receives a numerical **reward**
 - The goal is to optimize action-selection policy to maximise the total reward received
 - The learner needs to **explore** (try) the different actions and **exploit** the seemingly most rewarding arms
 - In practice, the learner has access to **contextual information** in each round to infer which action leads to the highest rewards

Contributions

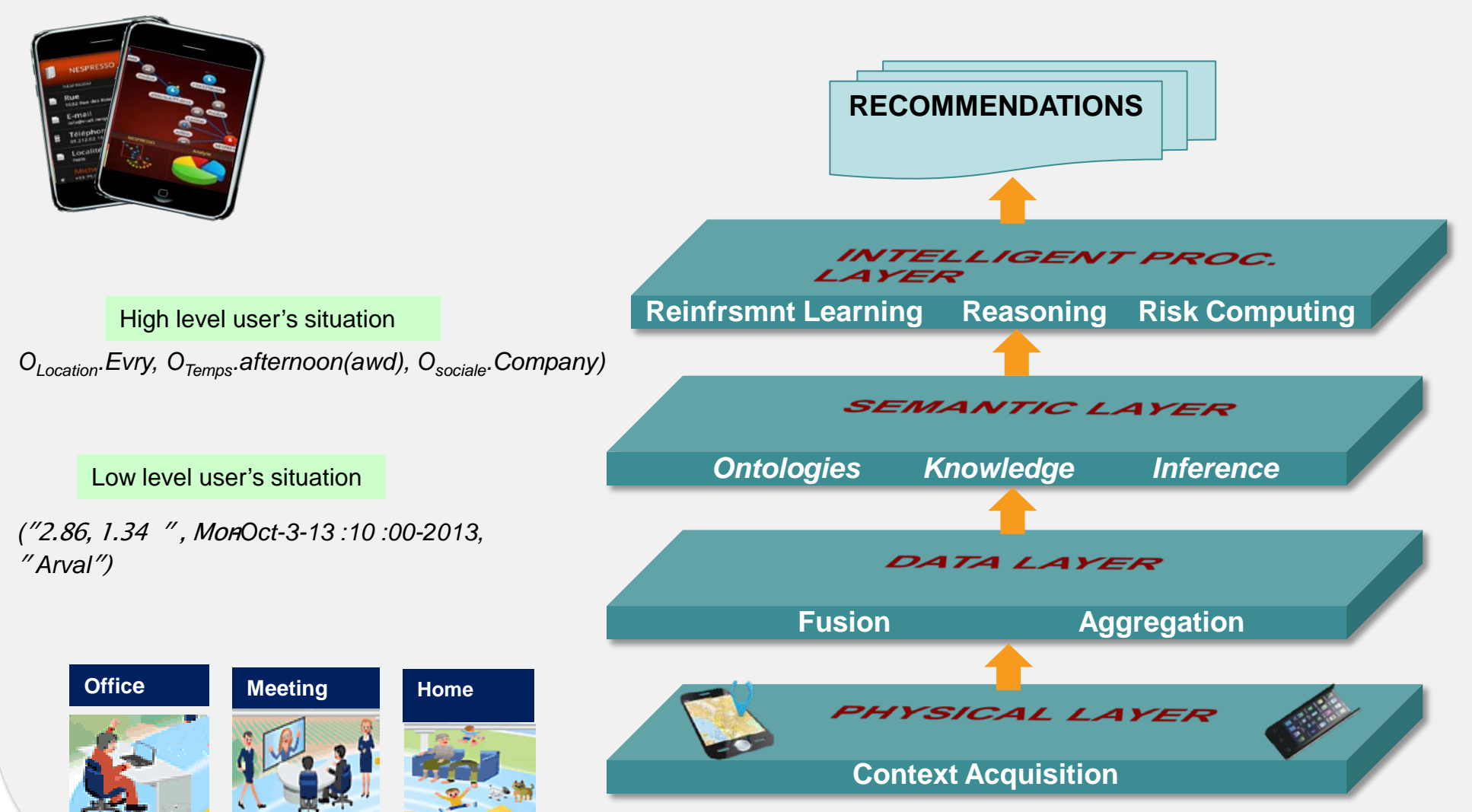
Models

- Modeling the context-aware recommender system as a bandit algorithm
- Modeling user, context, situation and risk

Algorithms

- A new semi-uniform strategy: contextual-epsilon-greedy strategy
 - Combining content-based filtering and reinforcement learning
- An algorithm R-UCB
 - Computes the probability of exploration by using the situation risk level $R(S)$
 - Three methods of risk computing
 - Using situations similarity (R_m)
 - Using situations concepts (R_c)
 - Using a Gaussian distribution (R_g)

General Approach



Evaluation Results

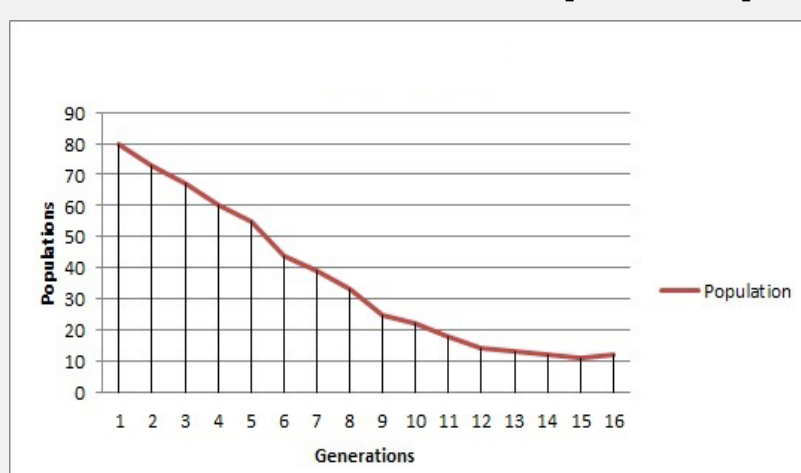
Data Set and Parameters

Nomalys Data Set

- 356 738 situations
- 5 518 566 navigations data
- 3500 users

Genetic Algorithm

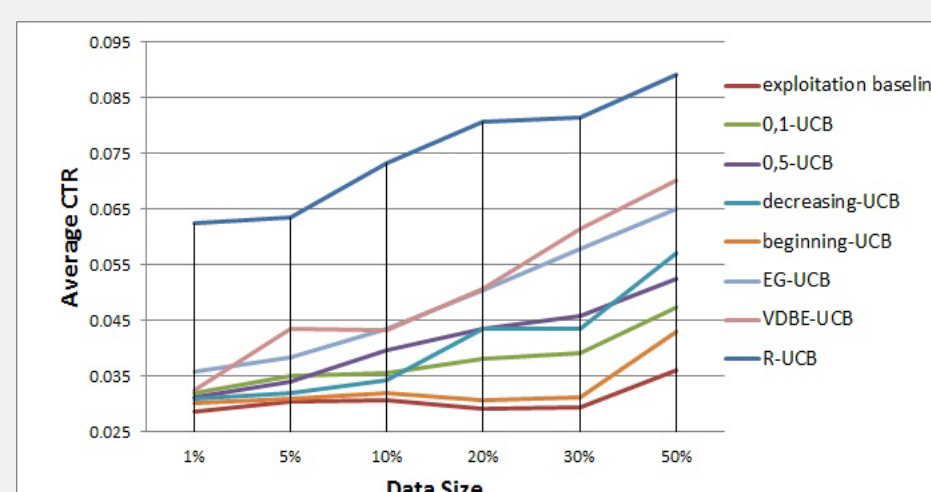
- Population: 80 Chromosomes
- Chromosomes: $\epsilon_{min}, \epsilon_{max}$, Threshold
- Results: $\epsilon_{min} \in [0.05; 0.13]$, $\epsilon_{max} \in [0.47; 0.56]$, Threshold $\in [0.7; 0.82]$



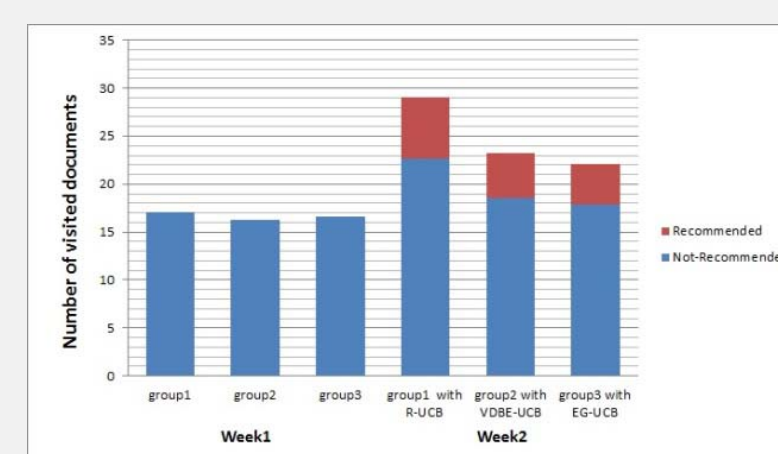
Offline Evaluation

Online Evaluation

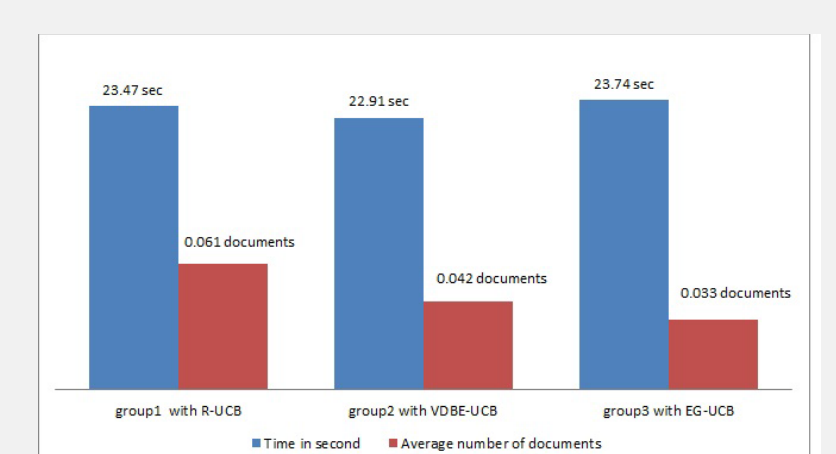
Data Size Variation



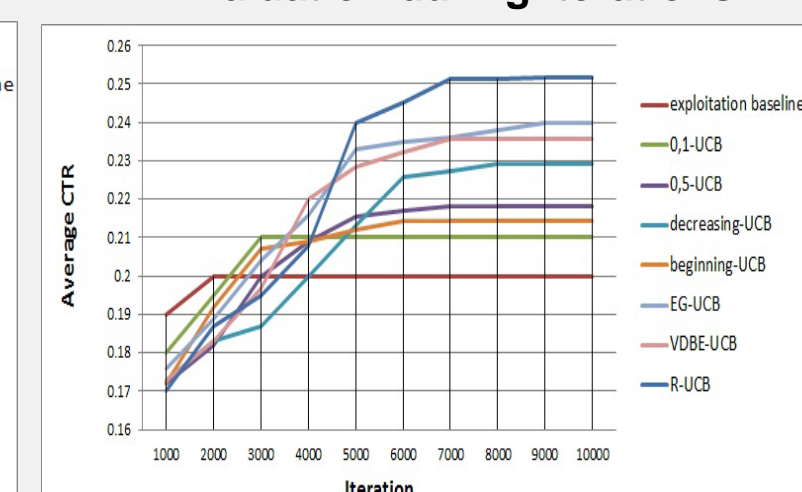
Number of visited documents



Time spent in documents



Evaluation during Iterations



Risk Variation

