

Mission Interdisciplinaire Défi MASTODONS

Journée ICUBE, 7 novembre 2014



AMADOUER : Analyse de MAsses de DOnnées de l'Urbain et l'Environnement

A. Baskurt – S. Servigne – JF. Boulicaut



UMR 5600



UMR 5008



EA 4126



UM5 5205

Contexte

■ MASTODONS

- « Big Data »
- « Data Science »



■ AMADOUER

- Absence de “commande” a priori
 - Masses de données hétérogènes et parfois volumineuses dans les domaines de l’urbain et de l’environnement
 - Questionnements variés
- Travailler à l’émergence de communautés « data science »
- Concepts – Généricité pour accompagner de bonnes ruptures méthodologiques en sciences

Singularités

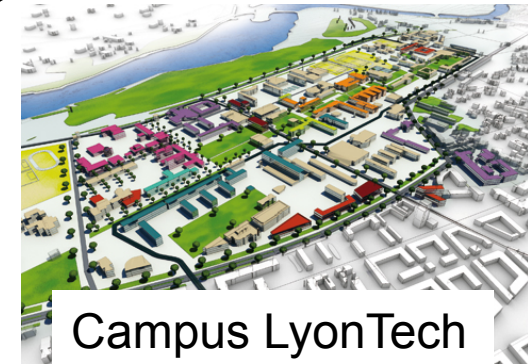
- Contexte local

MI

INEE
INSHS
INSIS
INS2I



GRAND LYON
communauté urbaine



CETHIL
UMR 5008



Bilan 2012

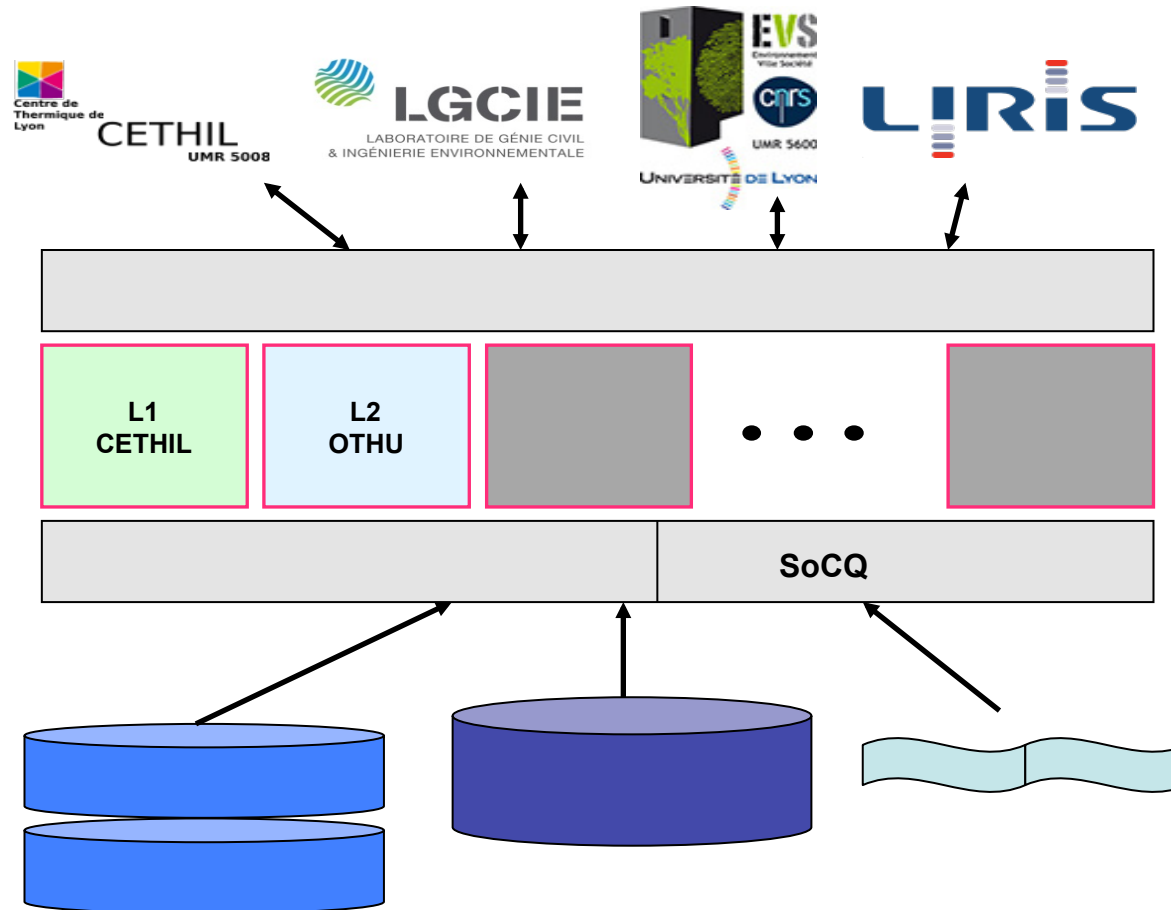
■ Identification “Données/Questionnement”

- EVS (INSHS, INEE)
 - Occulométrie dans un environnement équipé
- CETHIL (INSIS)
 - Combustion
 - Changements de phase
 - Monitoring d’installation de production d’énergie solaire photovoltaïque
 - Monitoring énergétique de bâtiment intelligent
- LGCIE-OTHU
 - Observatoire de Terrain en Hydrologie Urbaine
- LIRIS
 - Plate-forme SoCQ4Home (mi-lourd CNRS)
- GRAND LYON
 - Données : mobilité, énergie, cohésion sociale, etc

Flux vidéos

“Réseaux de capteurs”

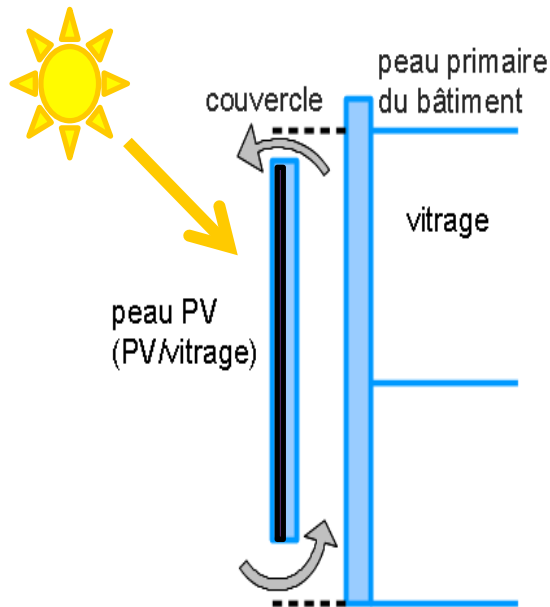
Plate-forme



Données CETHIL (projet Ressources)

Réseaux de capteurs (électricité, anémomètres, pyranomètres, thermocouples, station météo+suntracker) ... Toutes les 2 minutes sur 1

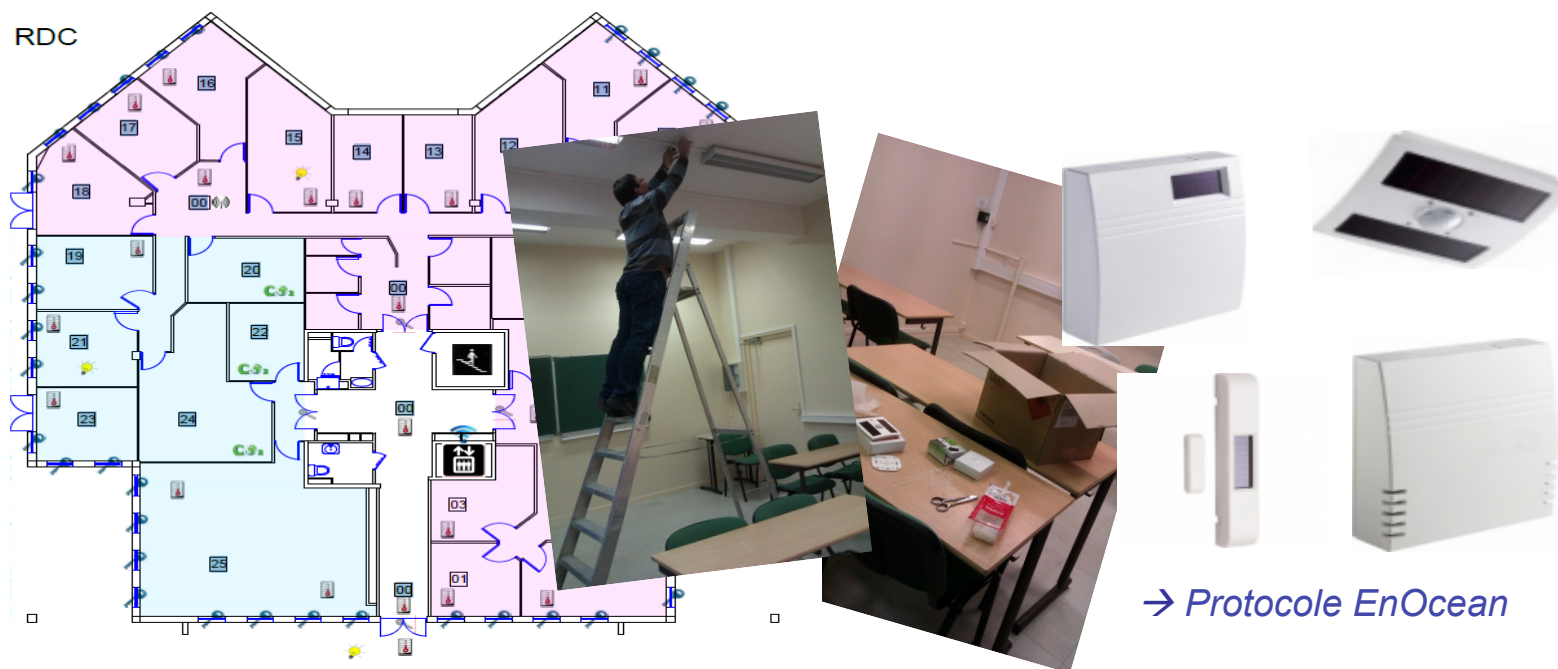
an



Questions prioritaires : Comportements “type” vs. “anormaux” (e.g., impact du vent), **loi empirique pour lier le débit d’air aux autres variables**

Données : Bâtiment intelligent

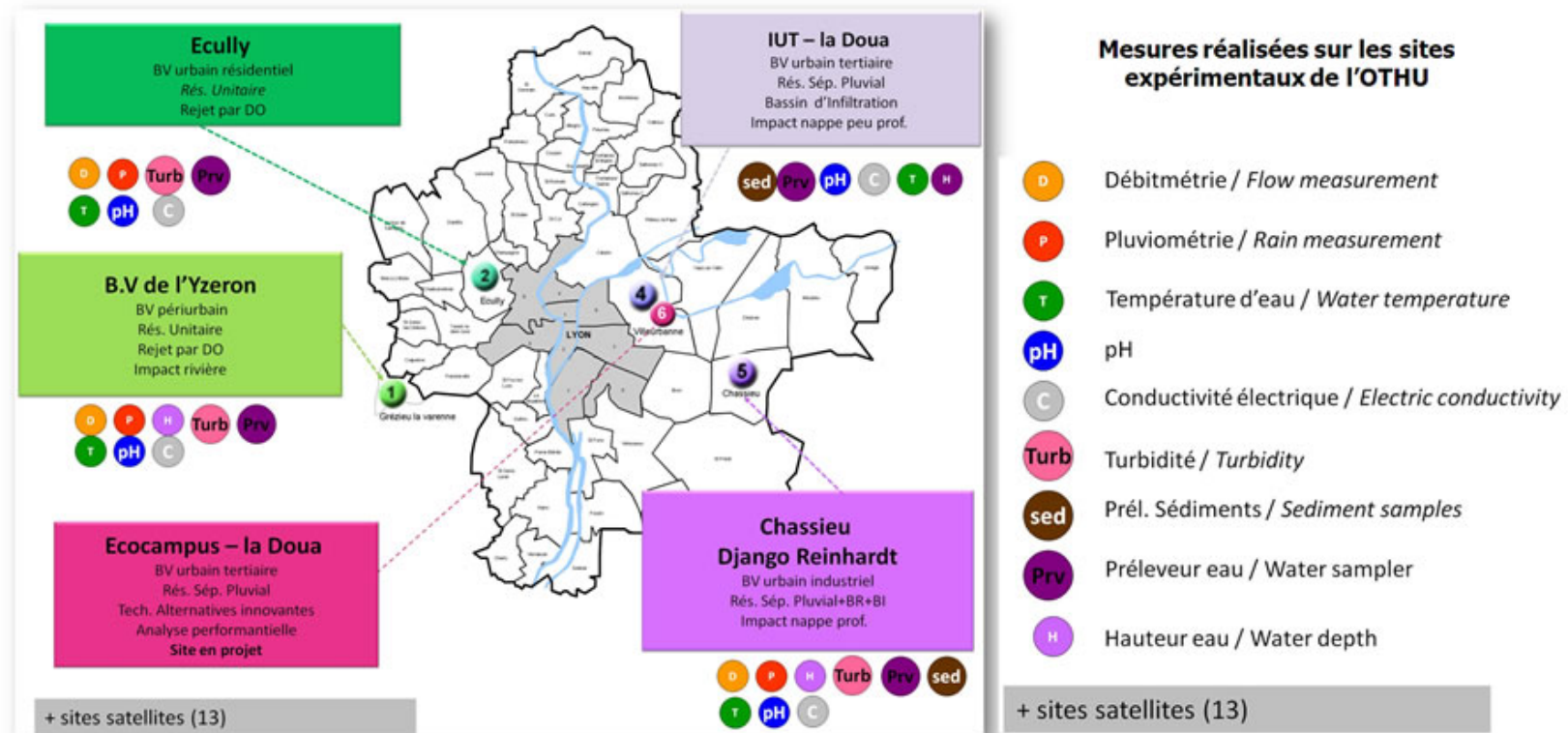
- Plateforme MARBRE (INSA : LIRIS-CETHIL-EVS)
 - Réseau 300 capteurs sans-fils (température, humidité, Co2, luminosité, ouverture porte/fenêtre)
 - Emission capteurs : toutes les 15 minutes ou selon événement
 - Enquêtes usagers



Question prioritaire : Validation de modèles physiques avec prise en compte des usages

Données OTHU

Réseaux de capteurs opérationnels depuis 2001 (Observatoire) :
toutes les 2 minutes 7j/7 24h/24



Questions prioritaires : qualité des données, “la constitution des flux polluants n’est pas comprise” (multi-échelles, nouvelles variables, ...)

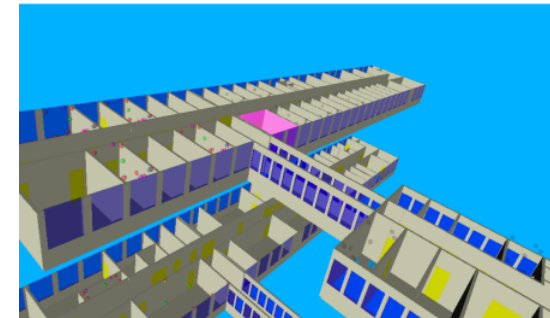
Verrous traités depuis 2013

- Modèles de données génériques adaptés aux flux
- Concepts et langage d'exploration données/flux
- Découverte de connaissances (fouille/modèles) : construction de (nouveaux) domaines de motifs
- Instanciation en thermique
 - Validation de modèles physiques
 - Détection d'anomalies
 - Interactions et corrélations indirectes
 - Conception de nouveaux modèles dynamiques (modélisations non-physiques, comportementales, etc.) pour une meilleure prédiction de la production instantanée ou cumulée (injection réseau, couplage bâtiment)

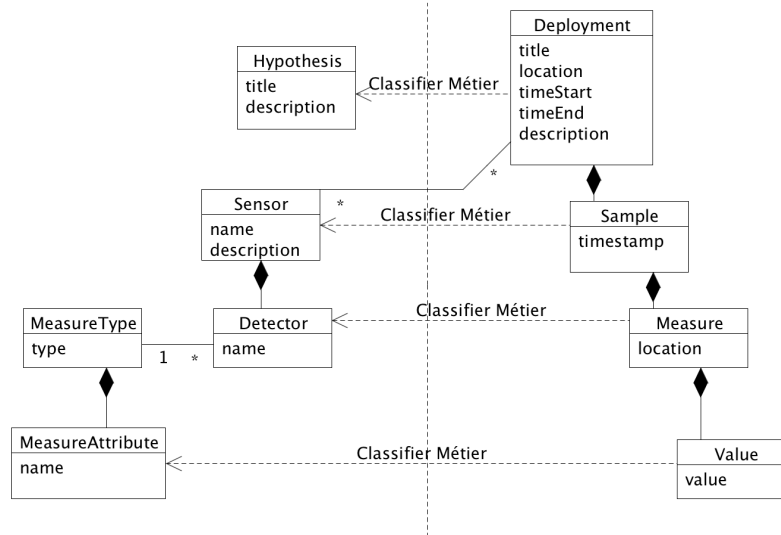
Résultats (1)

Modèles de données

- Définition et formalisation du modèle Virtual Generic Sensor (VGS) collaboration LIRIS (INS2I) et EVS (INSHS-INEE)
- Application à la prise en compte des usages en thermique du bâtiment.



Visualisation Web 3D



Sensor1	HumidityDetector1	10/02/2013 20:05:00	Humidity	Humidity	1000	(2,2,3)@502:309	42
Sensor1	HumidityDetector1	10/02/2013 20:05:01	Humidity	Humidity	1001	(2,2,3)@502:309	43
Sensor1	HumidityDetector1	10/02/2013 20:05:02	Humidity	Humidity	1002	(2,2,3)@502:309	44
Sensor1	HumidityDetector1	10/02/2013 20:05:03	Humidity	Humidity	1003	(2,2,3)@502:309	42
...
...
Sensor3	Question1	10/02/2013 20:05:00	Question	Humor	3000	(1,1,3)@502:411	5
Sensor3	Question1	10/02/2013 20:05:00	Question	Vitality	3001	(1,1,3)@502:411	7
Sensor3	Question1	10/02/2013 20:05:05	Question	Humor	3002	(1,1,3)@502:410	2
Sensor3	Question1	10/02/2013 20:05:05	Question	Vitality	3003	(1,1,3)@502:410	2

Résultats (2)

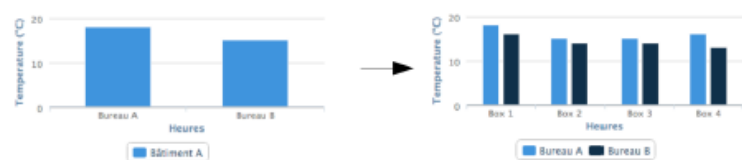
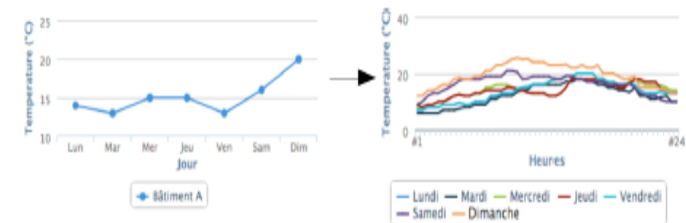
Exploration → Astéroïdes et Langage

- Définition et formalisation du concept d'Astéroïde
- Exploration multidimensionnelle interactive par agrégations dynamiques multi-échelles de données, pilotée par l'utilisateur

Schéma d'Astéroïde sur un schéma de Série Multidimensionnelle: sous-ensemble de dimensions, niveau "courant" de chaque dimension et schéma des valeurs

- Soit un Schéma de Série Multidimensionnelle $\mathcal{S} = (\mathcal{W}_S, \mathcal{V}_S)$
- Schéma d'Astéroïde $\mathcal{X} = (\mathcal{W}_X, [[l_j]]_{\mathcal{W}_X}^j, \mathcal{V})$ sur \mathcal{S}
 - $\mathcal{W}_X = \langle \mathcal{D}_{j_1}, \dots, \mathcal{D}_{j_m} \rangle \subseteq \mathcal{W}_S$ liste d'un sous-ensemble des dimensions du schéma de Série Multidimensionnelle $\mathcal{W}_S = \langle \mathcal{D}_1, \dots, \mathcal{D}_n \rangle$
 - $l_j \in L_{\mathcal{D}_j}$ niveau de la dimension $\mathcal{D}_j \in \mathcal{W}_X$
 - \mathcal{V} symbole de "relation" pour les valeurs des points
 - $schema(\mathcal{V}) \subset \mathcal{A}$

Note: un schéma de Série Multidimensionnelle $\mathcal{S} = (\mathcal{W}_S, \mathcal{V}_S)$ se "transpose" en schéma d'Astéroïde $\mathcal{X} = (\mathcal{W}_S, [[l_i^\bullet]]_{\mathcal{W}_S}^i, \mathcal{V}_S)$ avec l_i^\bullet le niveau de base de la dimension $\mathcal{D}_i \in \mathcal{W}_S$.



Project (i.e., select a subset of values and/or compute new values)

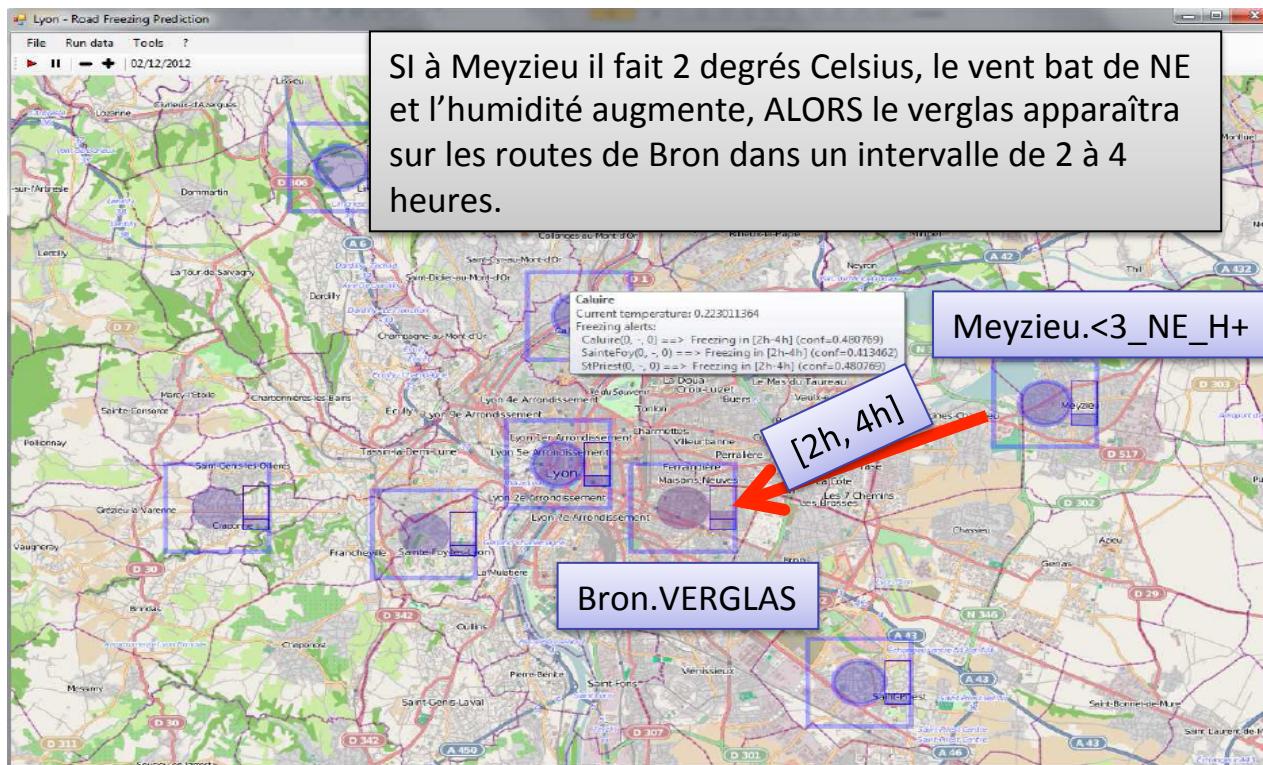
Calcul de nouvelles valeurs par l'application d'une fonction sur les valeurs de chaque point.

- Soit $x = ([[d_i]]_{\mathcal{W}_X}^i, [[M_i]]_{\mathcal{W}_X}^i, V)$ sur $\mathcal{X} = (\mathcal{W}_X, [[l_i]]_{\mathcal{W}_X}^i, \mathcal{V}_X)$
- Soit une fonction (simple) $\mathcal{F}_s^\circ : \mathcal{V}_X \rightarrow \mathcal{V}'$
- Résultat $x' = ([[d_i]]_{\mathcal{W}_X}^i, [[M_i]]_{\mathcal{W}_X}^i, V')$ sur $\mathcal{X}' = (\mathcal{W}_X, [[l_i]]_{\mathcal{W}_X}^i, \mathcal{V}')$
 - $V'(p) = \mathcal{F}_s^\circ(V(p))$
- Notation $x' = \Upsilon_{\mathcal{F}_s^\circ}(x)$ (*lettre v (upsilon) majuscule*)

Résultats (3)

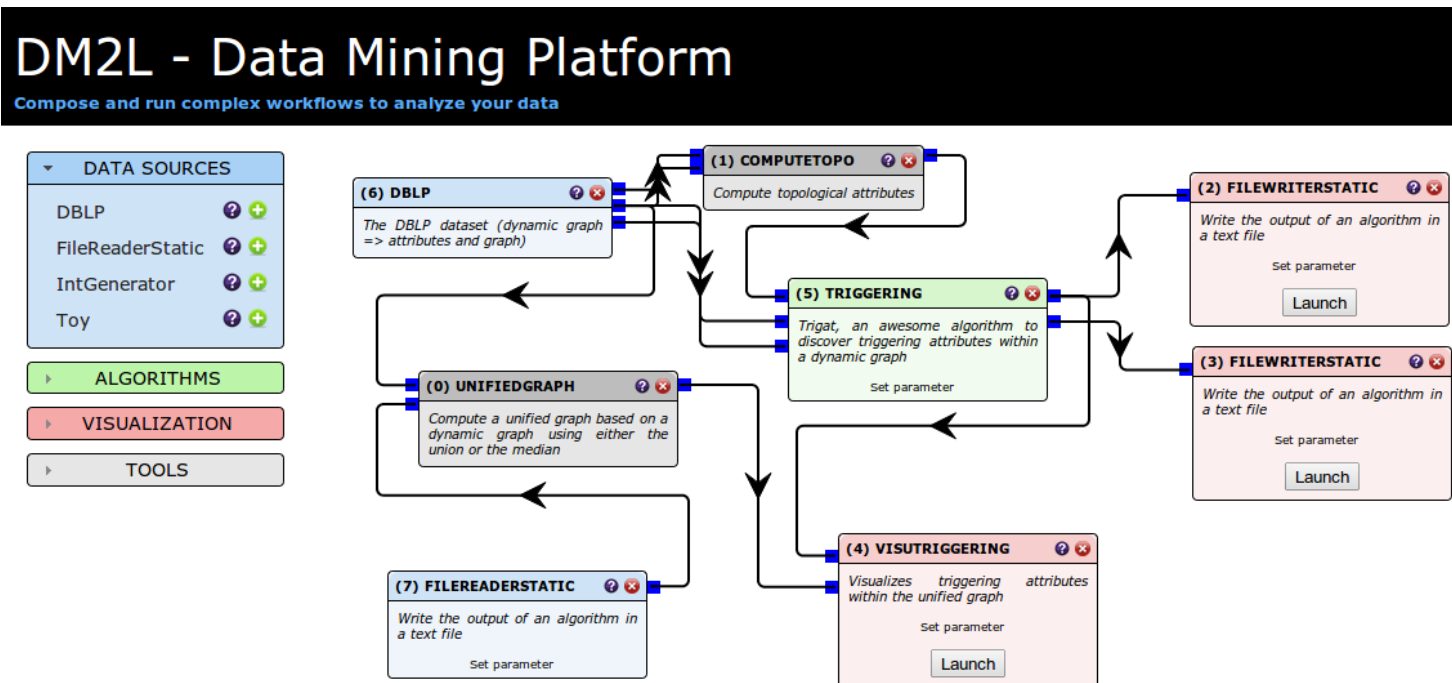
Extraction de connaissances → KDD

- Découverte de dépendances temporelles dans l'analyse de flux de données du Grand Lyon : application au salage de la voirie



Résultats (4)

Plateforme AMADOUER



A noter

- C'est dans la **généricité** des concepts proposés et implémentés que réside la valeur ajoutée
 - Langage de workflow pour une capitalisation des développements de prototypes en exploration et fouille des données
 - Langages et modèles adaptés aux données issues de réseaux de capteurs ... en évolution
 - Algorithmes de fouille de données génériques (e.g., analyse des dépendances temporelles entre flux de données)

Travail en cours sur 2014

- Développement de démonstrateurs sur la [plateforme AMADOUER](#)
- Synthèse sur nos analyses de masses de données urbaines produites par le [Grand Lyon](#) (bientôt métropole)
- Rapprochement avec [ANIMITEX](#) et travail sur une feuille de route concernant les verrous « [Hétérogénéité & Big Data](#) » (illustrations dans les sciences de l'environnement)
- [Coopération](#) avec le [PPME](#) de l'Université de Nouvelle Calédonie (Nazha Selmaoui, Frédéric Flouvat)
 - [Exploitation de modèles experts pour la dérivation de contraintes](#)

$$\{\tau \in L \mid q(\tau, r) \text{ est vrai}\}$$

Conclusion (1)

- Des expériences en cours sur l'émergence de “petites” communautés “Data Science”
 - Thermique, Energie, Environnement
 - Sciences du Vivant, Physique, Mathématique, Sciences Humaines et Sociales, etc
- A confronter à d'autres projets
 - Notamment dans le contexte MASTODONS
 - Actions avec [ANIMITEX](#) (Juin 2014, Novembre 2014)
- Tâche ingrate mais nécessaire ;-)

Conclusion (2)

- Il faut accompagner des ruptures méthodologiques de recherche en identifiant des solutions génériques
 - A l'échelle de, e.g., MASTODONS, il y a donc des **REX** dans des contextes très différents
 - La valorisation de masses de données identifiées (par exemple PETASKY) avec des défis très "Volume"
 - E.g., passage à l'échelle de tel ou tel algorithme sur les images LSST
 - L'identification des difficultés et la proposition d'éléments de méthodes pour **réussir des expériences en science des données**.
 - E.g., méthodologie de co-construction et la validation de domaines de motifs

Co-construction de domaines de motifs

- Le propriétaire des données au coeur des processus de création de la valeur ajoutée

(L,C,M,r)

Langage de motifs

Contraintes Primitives

Modèle des connaissances a priori

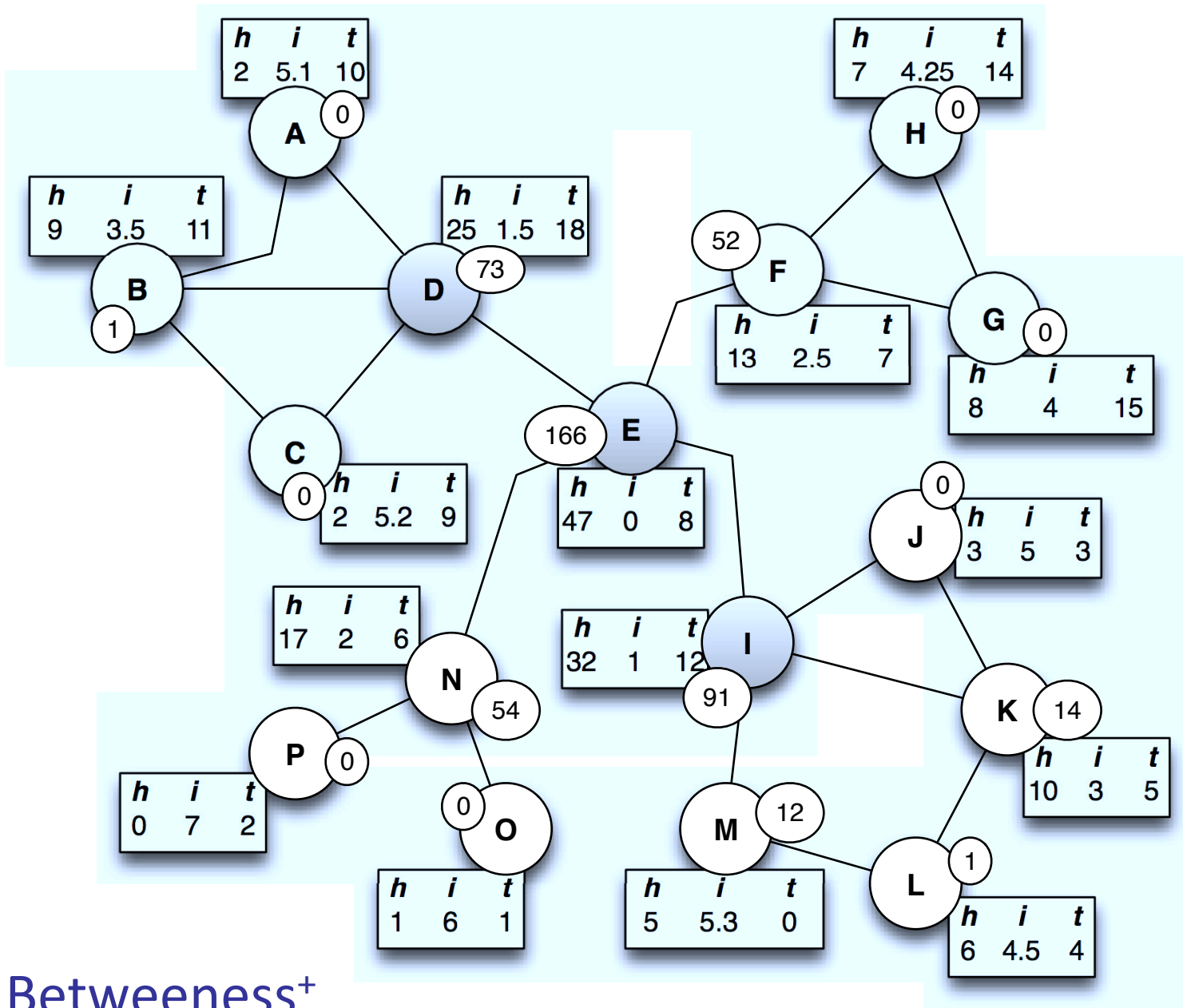
Données



Requêtes inductives $\{\tau \in L \mid q(\tau, M, r) \text{ est vrai}\}$

- Comment concevoir et prototyper un domaine de motif ?

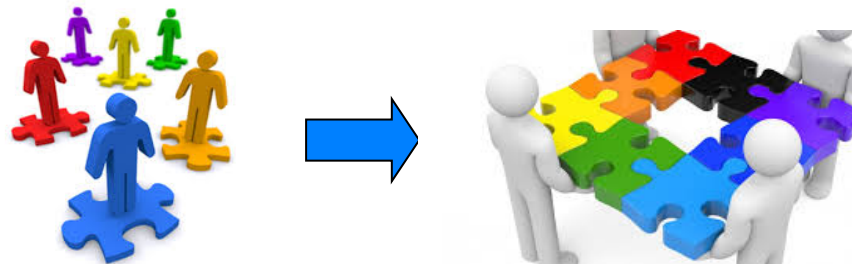
Déclaratif vs. Calcul, Générique vs. Ad-hoc, Exact vs. Approché, Préférences et souplesse, etc



$h^+ i^-$ Betweenness⁺

Conclusion (3)

- Valorisation de notre REX sur l'émergence de communautés en science des données
 - Impact sur les formations ?
 - GDR ?



Participants AMADOUER



M. Cottet, J-M. Deleuil, J-Y. Toussaint, S. Vareilles



J. Bonjour, L. Gaillard, S. Giroux, J. Jay, C. Ménézo, H. Pabiou



A. Baskurt, J-F. Boulicaut, G. Brochet, S. Fenet, Y. Gripay, M. Keytoue, M. Plantevit, Y. Pitarch, M. Scuturici, C. Robardet, S. Servigne, C. Wolf (PPME N. Selmaoui, F. Flouvat – Université de Nouvelle Calédonie)



S. Barraud, N. Walcker, L. Bacot