

Experimental Evaluation of a Dynamic Cubing system: workflow, metrics and prototype

Anne Tchounikine, Maryvonne Miquel, Usman Ahmed

LIRIS

CNRS UMR 5205, INSA-Université de Lyon, France

Motivations

- Motivé par des travaux précédents : Cube et OLAP dynamique
- Problématique de l'évaluation de la solution
 - Démonstration et qualification formelle
 - Démarche expérimentale
 - Prototypage
 - Appréciation des résultats
- **Démarche expérimentale**
 - Observer des comportements
 - Vérifier certaines intuitions
 - Mesurer des résultats

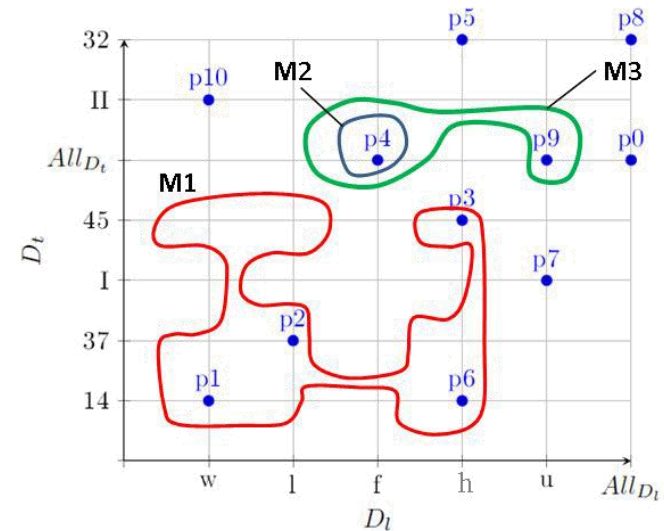
Objectifs et éléments de la contribution

- Résultats et **démarche** de l'évaluation expérimentale
- Un prototype et des expérimentations pour :
 - Montrer la faisabilité
 - Effectuer des tests fonctionnels
 - Ajuster les paramètres de la solution
 - Evaluer la performance
 - Etudier le comportement
 - Mener des études comparatives
- Définition d'un workflow pour l'expérimentation
 - Métriques
 - Caractérisation des éléments en entrée
 - Définition de scénarios d'exécution
 - Outils de monitoring

→ Appliqué à notre proposition de cube dynamique

Le modèle de cube

- Contexte : BI temps-réel
- Un espace multidimensionnel hiérarchique
 - Tous les membres d'une dimension sur le même axe
 - Membres non ordonnés
 - Axes construits dynamiquement
- Les **MBS** (Minimum Bounding Set) regroupent des points de même niveau hiérarchique



p4 \odot p2

M1 = {45, 37, 14} \times {w, l, h}

|M1| = 9

||M1|| = 4

Let $\Delta = \{(x_{1,1}, x_{1,2}, \dots, x_{1,n}), (x_{2,1}, x_{2,2}, \dots, x_{2,n}), \dots, (x_{m,1}, x_{m,2}, \dots, x_{m,n})\}$ be a set of m points lying in a same n -dimensional hyper-plane i.e. $x_{j,i} \in \text{domain}(l_i^{k_i})$ for $1 \leq i \leq n$, $1 \leq j \leq m$ and k_i is a level in the hierarchy of dimension D_i . A minimum bounding space (MBS) constructed over Δ , denoted by M_Δ , is defined as:

$$M_\Delta = \bigcup_{j=1}^m x_{j,1} \times \bigcup_{j=1}^m x_{j,2} \times \dots \times \bigcup_{j=1}^m x_{j,n}$$

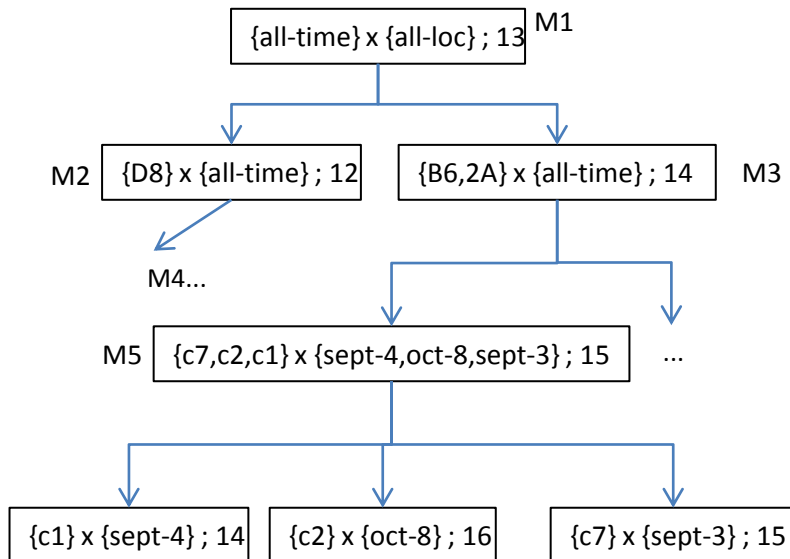
where each set $\bigcup_{j=1}^m x_{j,i}$ is called the i^{th} ($1 \leq i \leq n$) dimension edge E_i of M_Δ

$|M_\Delta| = \prod_{i=1}^n |E_i|$ the number of points covered by the MBS M_Δ ,

$||M_\Delta|| = |\Delta|$ the number of points enclosed in the MBS M_Δ .

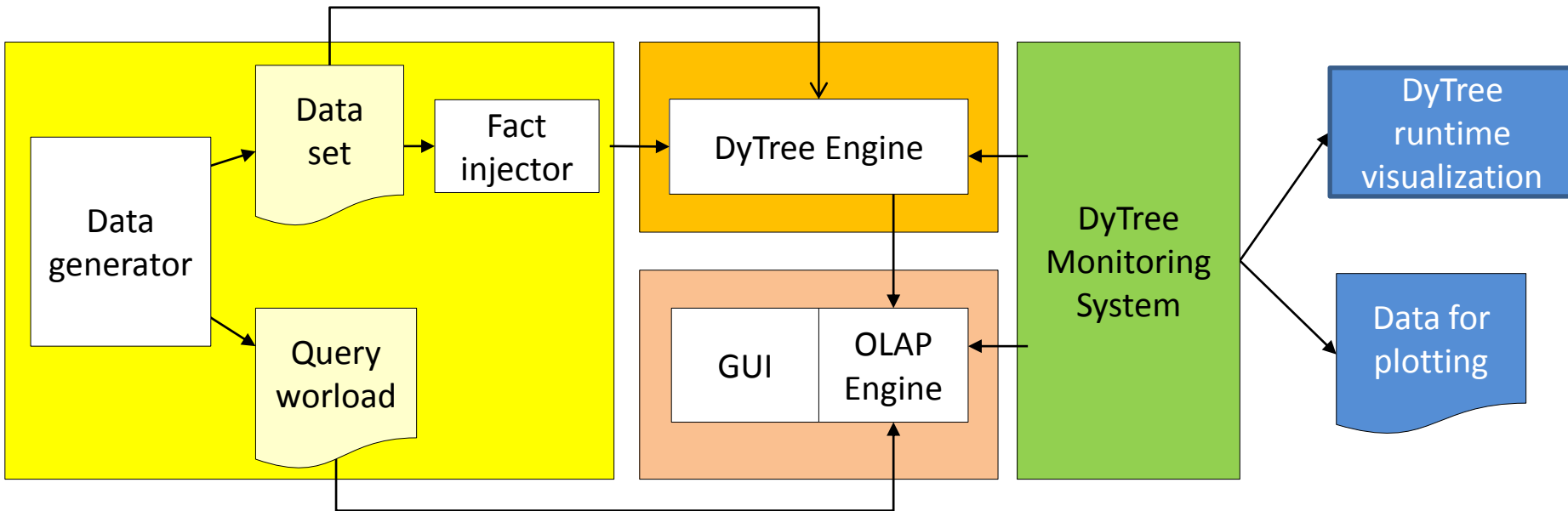
Le modèle d'arbre

- Le **DyTree** est le cube matérialisé
- Les points (faits) sont positionnés dans l'espace au fur et à mesure de leur arrivée. Ils sont les feuilles de l'arbre
- Les points sont regroupés dans un MBS. Lorsqu'un MBS est plein :
 - il est divisé (split)
 - un MBS de points couvrants est créé
 - La mesure agrégée est calculée



→ les MBS sont des morceaux de cuboïdes
→ les MBS sont de plus en plus détaillés
→ les MBS créés sont naturellement denses

Le prototype



- Efficience du Dytree et des algorithmes associés ?
- Comportement en fonction des contextes d'exécution ?

Search Query: "Find Sum (Quantity)"

Search Results

Dimensions:

- Customer
 - Supplier Dimension
 - Region
 - Country
 - City
 - Supplier
- Part Dimension
 - Category
 - Brand
 - Part
- Time Dimension
 - Year
 - Month
 - Day
 - Hour

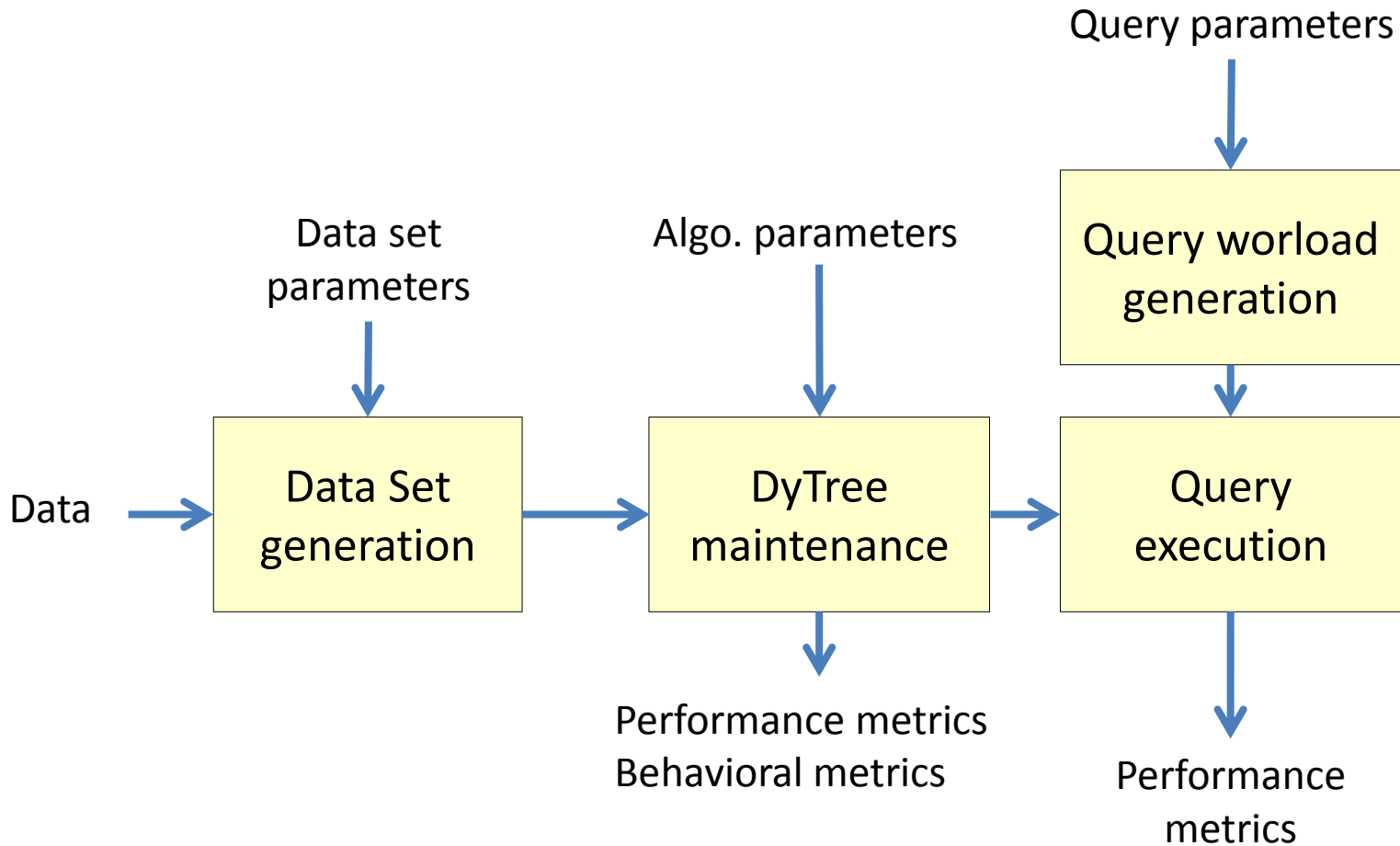
Rows: Customer Region, Supplier Country

Columns: Time Month, Part Category

		January, 1999				
		MFGR11	MFGR25	MFGR31	MFGR32	MFGR33
Europe		12520	4875	22020	5815	275
America		5865	55340	15940	44970	151
	France	0	0	0	0	0
	Jamaica	0	1265	0	350	0
	Tunisia	0	1605	695	185	250
	Algeria	35	5820	6350	1030	325
	Nigeria	2740	2125	4280	750	185
	South Africa	4060	160	550	265	210
	Kenya	0	545	0	405	0
	Saudia Arabia	40	225	140	175	155
	UAE	420	565	1105	855	300
	Iraq	820	1150	905	535	40
	Kuwait	255	345	675	395	235
	Germany	0	0	0	0	0
Africa	Oman	90	300	580	315	550
	Dubai	60	3615	840	1860	315

Query Execution Time: 238340 Mill-Seconds

Workflow de l'expérimentation



Métriques de performance

- Métriques de performance
 - Temps d'exécution
 - Temps cumulé pour l'insertion d'un ensemble de faits
 - Temps d'insertion atomique d'un fait
 - Temps de réponse aux requêtes
 - Taille de l'arbre en mémoire
 - Les feuilles (faits) sont stockés sur disque
 - Les nœuds (agrégats) en mémoire

Métriques de comportement

- Métriques sur l'arbre
 - Nombre de nœuds internes
 - Nombre de super nœuds
 - Profondeur
 - Largeur
- Métriques sur les nœuds
 - Taux de remplissage

$$\textit{fillRatio}(\textit{node}) = \frac{\textit{number of children}(\textit{node})}{\textit{directory node capacity}}$$

- Volume et Densité
 - Calculé sur le résultat (N) d'un drill-down des coordonnées du MBS (M)
 - Nombre de points couverts et nombre de points inclus

$$\textit{volume}(M) = |N| \text{ and } \textit{density}(\textit{node}) = \frac{\|N\|}{\textit{volume}(M)}$$

Métriques vers indicateurs

- **Transformer les mesures en indicateurs (→ signification, impact)**

- Métriques de performance

- Temps d'exécution

- Insertion d'un fait, temps cumulé, requête

- Taille de l'arbre en mémoire

→ efficacité de la solution dans un contexte BI "agile"

- Métriques de comportement

- Métriques portant sur l'arbre

- Nombre de nœuds internes
- Nombre de super nœuds
- Profondeur
- Largeur

→ nombre d'agrégats matérialisés

- Métriques portant sur les nœuds

- Taux de remplissage
- Volume
- Densité

→ granularité des agrégats

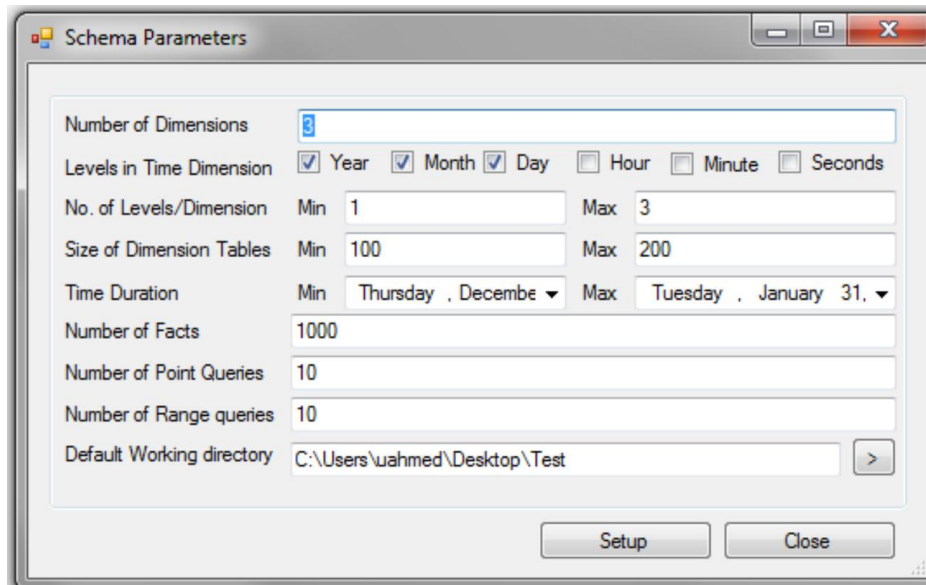
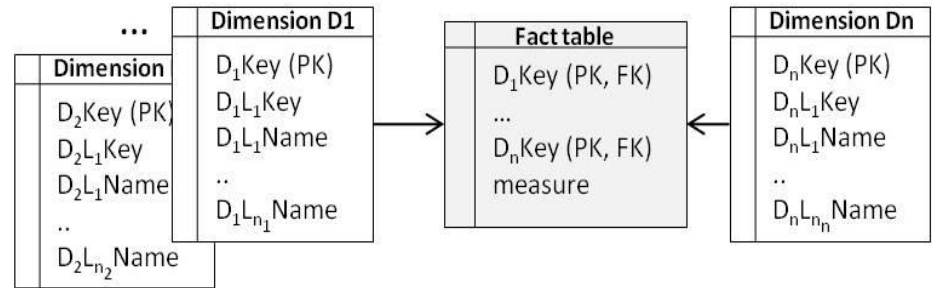
→ pertinence des agrégats

Éléments en entrée

- Lister et caractériser les inputs pour créer des scénarios d'exécution
- Paramétrage des algorithmes
 - Capacité des nœuds
 - Recouvrement (overlap)
- Paramétrage des data sets
 - Schéma
 - Taille (cardinalité des dimensions, volume de faits)
 - Densité
 - Ordre
 - Fréquence d'arrivée
- Paramétrage des requêtes
 - Nombre de requêtes
 - Types (point, range, group by)

Schéma des data sets

- Benchmark : Star Schema Benchmarks
- Données synthétiques
- Données réelles : SoQ4Home
 - Bâtiments intelligents
 - Capteurs de température
 - Déployé sur le campus LyonTech

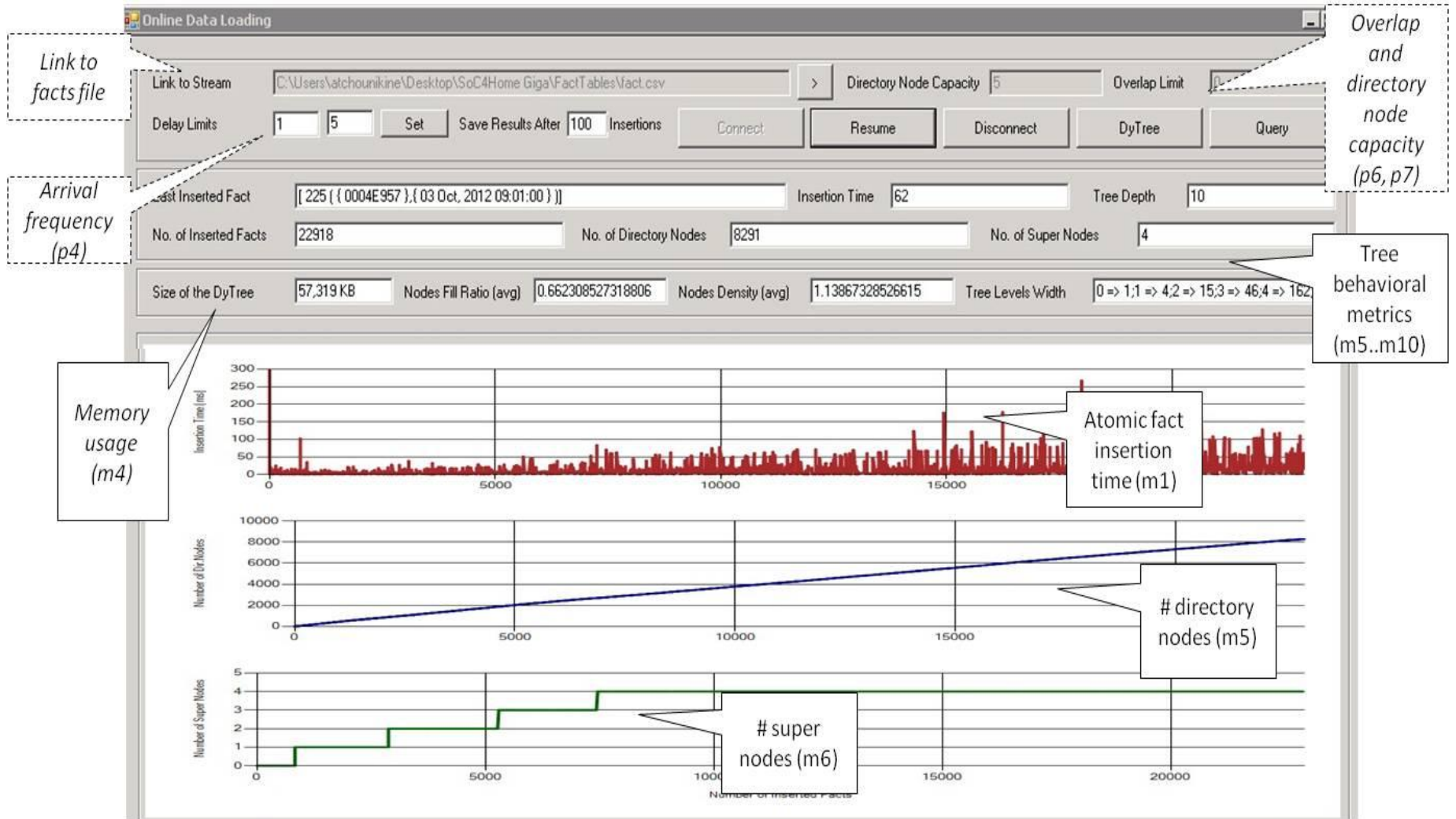


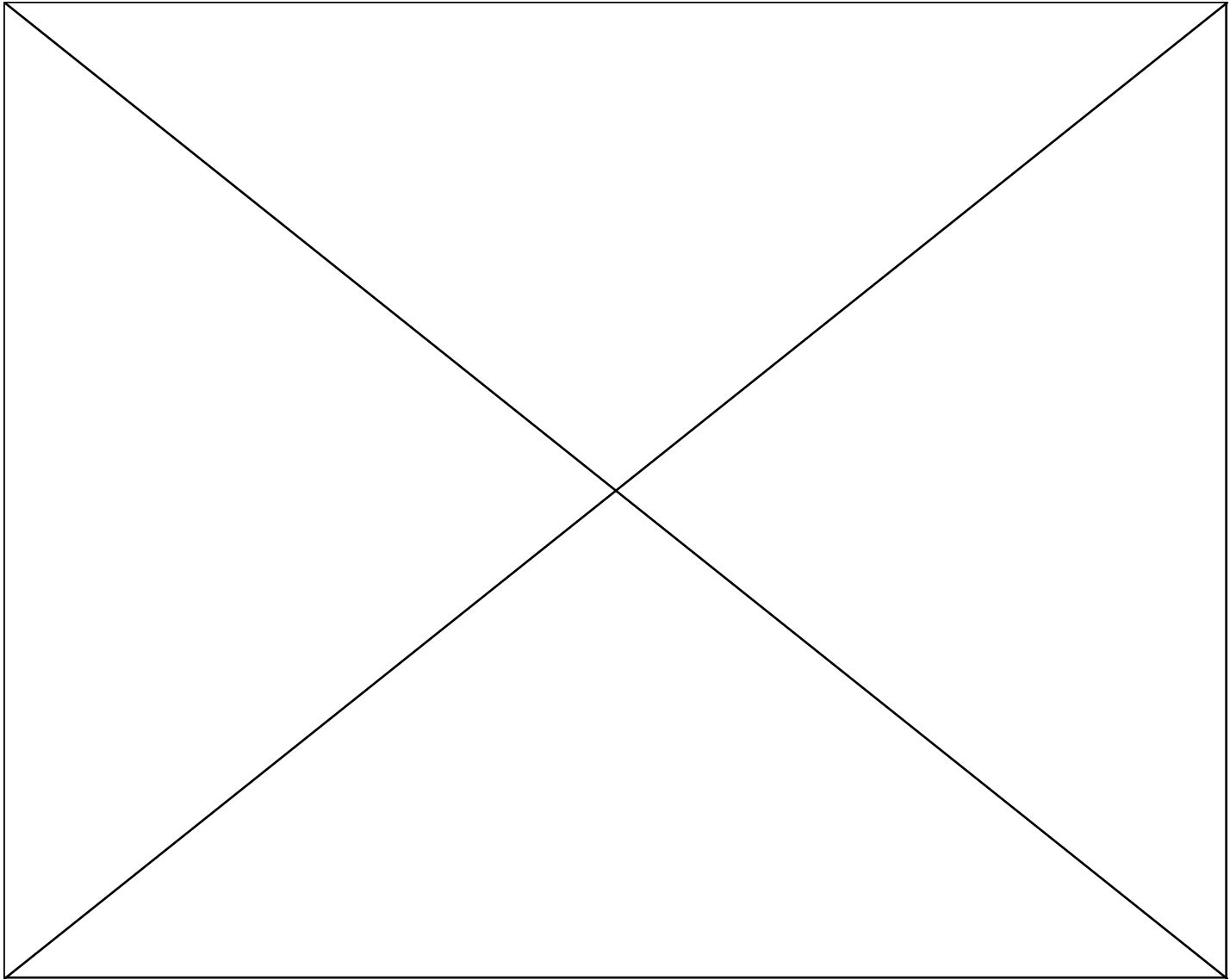
Constitutions de scenarios

- En faisant varier les paramètres en entrée

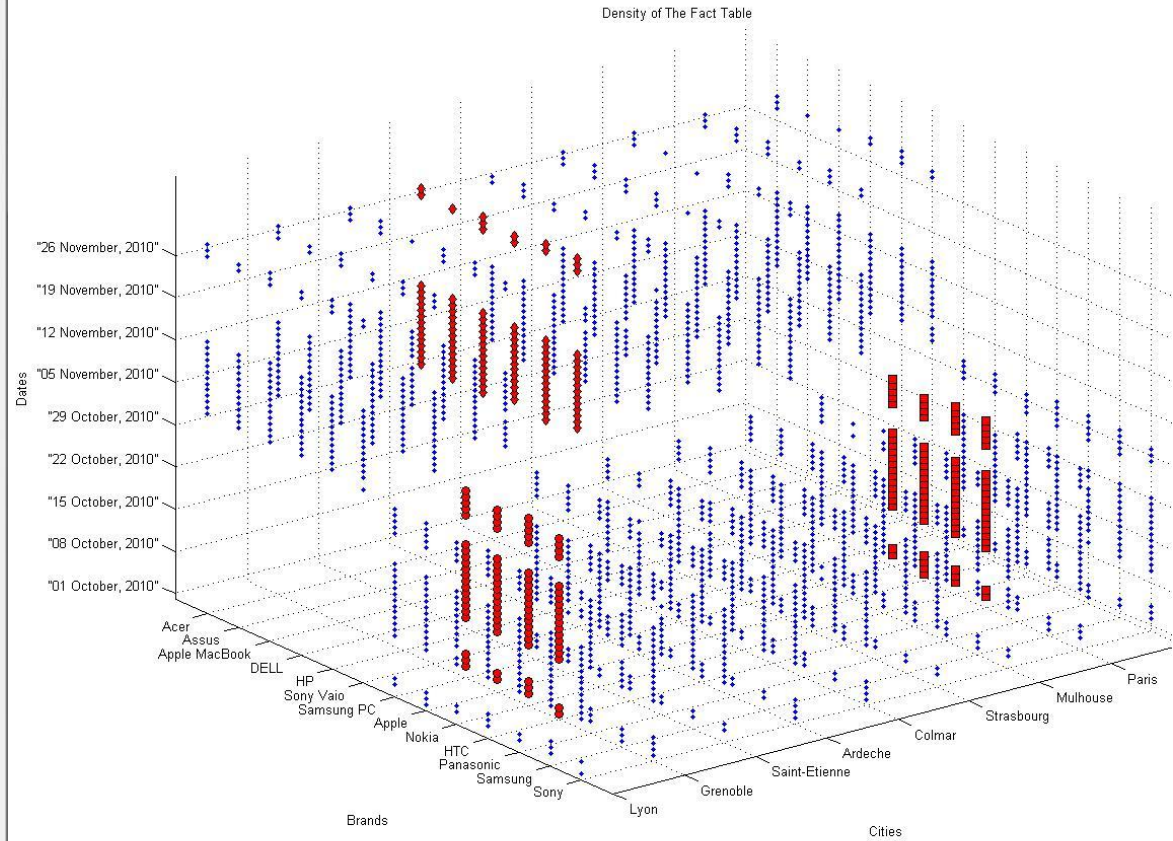
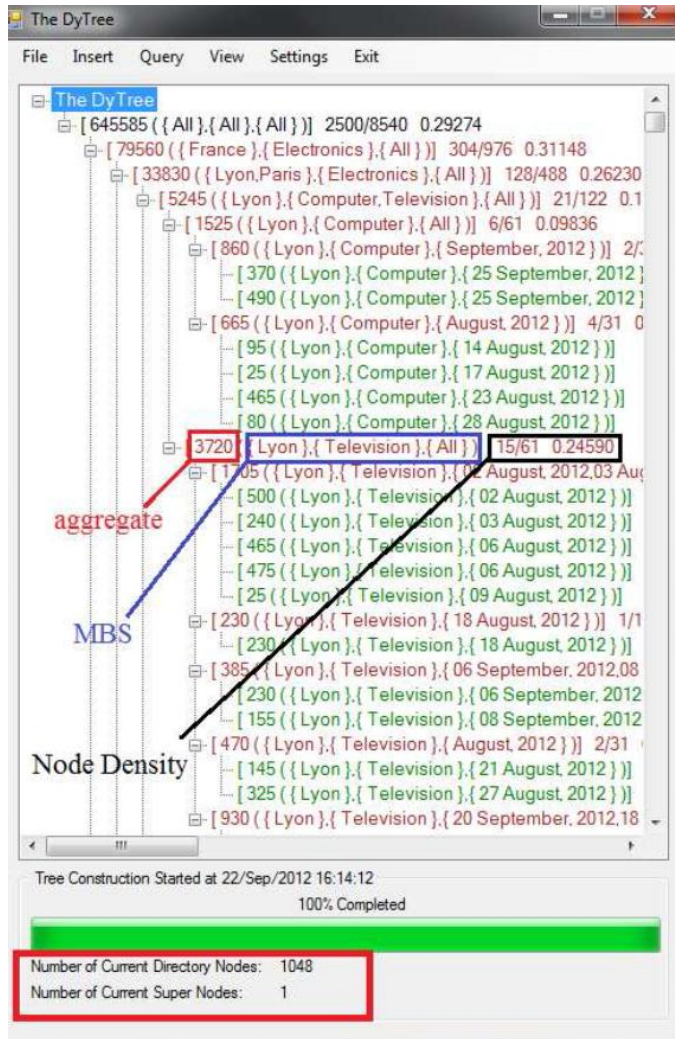
Scenario	Data Set	# dim	# fact	density	arrival freq	order	dir cap	overlap
Varying overlap	SSB	4	10M	calc.	na	na	15	0 to 15
Varying dir. cap.	SSB	4	10M	calc.	na	na	5 to 75	0
Scaling # facts	SSB	4	up to 10M	calc.	na	na	15	5
	synth	10	up to 100M	calc.	na	0%	15	5
	Soq4Home	2	na	na	na	na	15	5
Scaling # dim.	synth	2 to 30	up to 100M	calc.	1 to 5 ms	0%	15	5
Varying density	synth	10	calc.	0.2 to 0.6	1 to 5 ms	0%	15	5
Delayed arrival	synth	10	up to 100M	calc.	1 to 5 ms	5 to 80%	15	5

Outils de monitoring : en temps réel

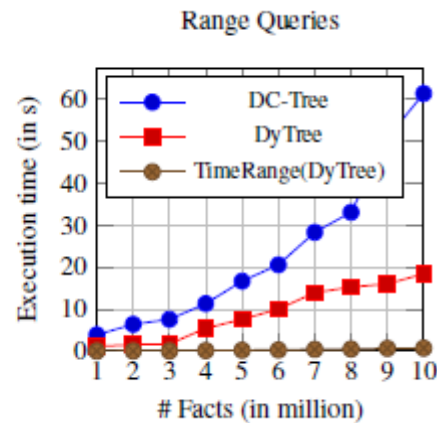
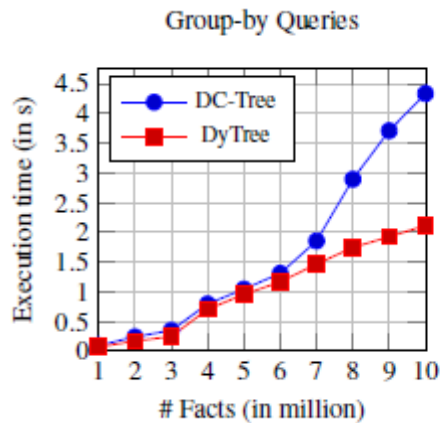
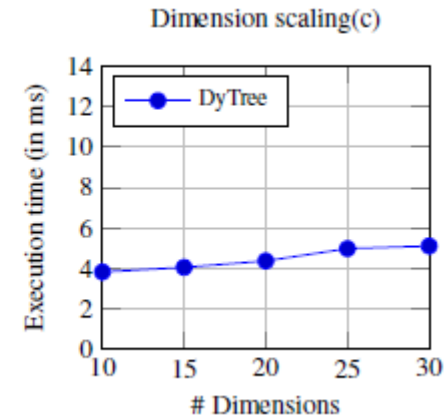
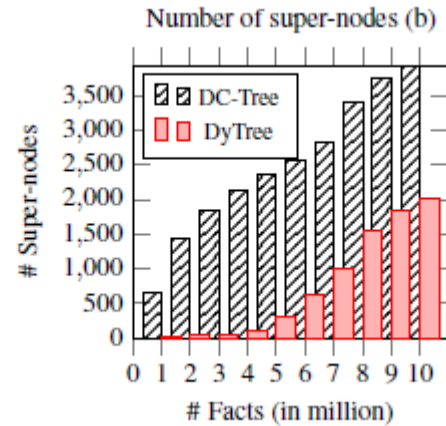
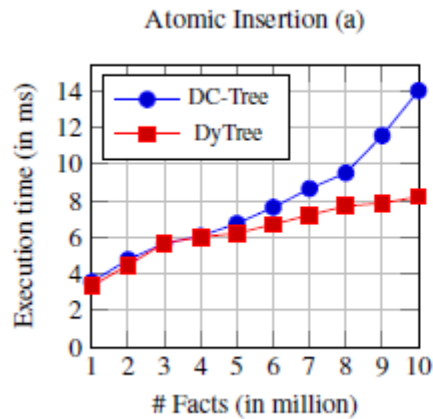




Outils de monitoring : des instantanés



Outils de monitoring : analyses a posteriori



Conclusions

- Ce prototypage nous a permis d'illustrer le modèle
- Une démarche expérimentale permet de valoriser le prototypage
 - Instancié sur notre exemple mais adaptable
- Approche complémentaire aux approches benchmark
- Perspectives
 - Les requêtes
 - Type de requêtes
 - Volume des résultats...
 - **Data set synthétiques**
 - **Génération de faits sous contraintes**