Raisonner avec la provenance sur les données du Web

Antoine Amarilli

Concours CNRS 06/03

6 avril 2016

Parcours

2013-2016: Thèse à Télécom ParisTech avec Pierre Senellart :

- · Tirer parti de la structure des données incertaines
- · Soutenue le 14 mars 2016

2012-2013: Pré-doc:

- · 3 mois à **Tel Aviv** avec Tova Milo
- 5 mois à Oxford avec Michael Benedikt

2009-2013 : École normale supérieure de Paris, master MPRI

Publications

Bases de données :

- · ICDT'14 (prédoc à Tel Aviv)
- · PODS'16 (thèse)

Logique et automates :

- · ICALP'15 (thèse)
- · LICS'15 (thèse)

Intelligence artificielle:

- IJCAI'15 (pré-doc à Oxford)
- IJCAI'16 (avec Oxford)
- 7 autres publications internationales avec comité de lecture

• 1 brevet avec Google New York (stage de M1)

Résumé des travaux antérieurs

Interrogation de données relationnelles incertaines

Vue d'ensemble : Données relationnelles incertaines

Évaluer une requête logique sur une base de données relationnelle

Problème : On ne dispose pas toujours des données exactes :

- · Données créées par des méthodes faillibles et non-exhaustives
- · Données annotées par des techniques d'apprentissage
- · Données bruitées ou périmées

→ Gérer les données relationnelles avec leur incertitude

Problème: Tâche souvent infaisable voire indécidable

• 1. données incomplètes • 2. données probabilistes

Données:

jour	type
9	congés
10	réunion
11	congés
18	congés
18	réunion

Requête logique:

Quelles réunions sont pendant mes congés?

Données:

jour	type	
9	congés	_
10	réunion	
11	congés	
18	congés	
18	réunion	
		-

Requête logique: ——

Quelles réunions sont pendant mes congés?

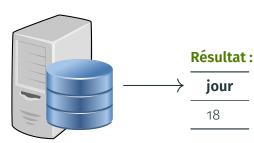


Données:

		_
jour	type	_
9	congés	_
10	réunion	
11	congés	
18	congés	
18	réunion	
		- ~

Requête logique : ——

Quelles réunions sont pendant mes congés?



Données:

jour	type	
9	congés	_
10	réunion	
11	congés	
18	congés	
18	réunion	
		-

Requête logique: ——

Quelles réunions sont pendant mes congés?



Données:

jour	type	
9	congés	
10	réunion	
11	congés	
18	congés	
18	réunion	
		-

Requête logique: -

Quelles réunions sont pendant mes congés?

Règles logiques :

 Je ne reviens pas pour un seul jour



Données:

		_
jour	type	
9	congés	_
10	réunion	
11	congés	
18	congés	
18	réunion	

jour type	
9 congés	
10 réunion	l
11 congés	
18 congés	
18 réunion	ı

Requête logique :

Quelles réunions sont pendant mes congés?

Règles logiques :

 Je ne reviens pas pour un seul jour



Données:

		_
jour	type	
9	congés	_
10	réunion	
11	congés	
18	congés	
18	réunion	

jour	type
9	congés
10	réunion
11	congés
18	congés
18	réunion
10	congés

Requête logique :

Quelles réunions sont pendant mes congés?

Règles logiques :

 Je ne reviens pas pour un seul jour



Données:

jour	type	-
9	congés	-
10	réunion	_
11	congés	
18	congés	
18	réunion	_

jour type

9 congés
10 réunion
11 congés
18 congés
18 réunion
10 congés

Requête logique :

Quelles réunions sont pendant mes congés?

Règles logiques :

 Je ne reviens pas pour un seul jour



Résultat:

jour

18

Données:

jour	type	
9	congés	-
10	réunion	_
11	congés	
18	congés	
18	réunion	_ ```

jour type

9 congés
10 réunion
11 congés
18 congés
18 réunion
10 congés

Requête logique :

Quelles réunions sont pendant mes congés?

Règles logiques :

 Je ne reviens pas pour un seul jour



Résultat:

jour

18 **10**

Résumé : raisonner sur les données incomplètes

→ Problème fondamental en intelligence artificielle :

Quelles réponses à la **requête** de l'utilisateur sont vraies dans toutes les complétions des **données** qui satisfont des **règles logiques**?

Approches existantes : Langages de règles décidables en IA :

- · Uniquement sur des graphes de données
- · Autorisant des complétions infinies

Résumé : raisonner sur les données incomplètes

→ Problème fondamental en intelligence artificielle :

Quelles réponses à la **requête** de l'utilisateur sont vraies dans toutes les complétions des **données** qui satisfont des **règles logiques**?

Approches existantes : Langages de règles décidables en IA :

- · Uniquement sur des graphes de données
- · Autorisant des complétions infinies
- → J'ai transposé ces résultats aux bases de données :
 - Étendre aux hypergraphes [Amarilli, Benedikt, IJCAI'15]
 - · Restreindre aux complétions finies [Amarilli, Benedikt, LICS'15]

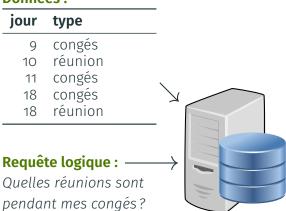
Données:

jour	type
9	congés
10	réunion
11	congés
18	congés
18	réunion

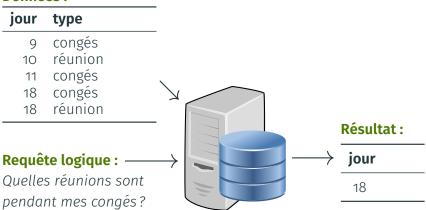
Requête logique:

Quelles réunions sont pendant mes congés?

Données:



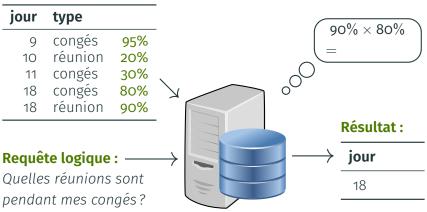
Données:



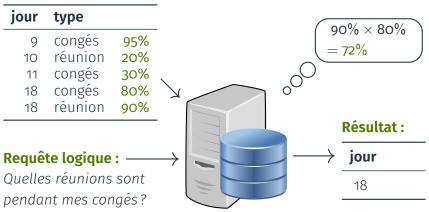
Données:

jour	type			
9	congés	95%		
10	réunion	20%		
11	congés	30%		
18	congés	80%	7	
18	réunion	90%		
				Résultat :
Requê	te logique	: —	\rightarrow	\longrightarrow jour
Quelle	s réunions	sont		18
pendai	nt mes cor	igés?		

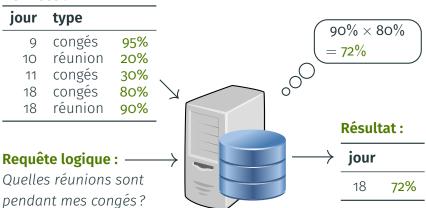
Données:



Données:



Données:



→ Problème d'évaluation de requêtes sur données **probabilistes** :

Quelle est la **probabilité totale** de chaque réponse quand les faits sont présents ou absents **indépendamment** avec la **probabilité** indiquée?

→ Problème d'évaluation de requêtes sur données **probabilistes** :

Quelle est la **probabilité totale** de chaque réponse quand les faits sont présents ou absents **indépendamment** avec la **probabilité** indiquée?

Approche existante (intensionnelle):

Calendrier						
id	jour	type				
t ₁	9	congés				
t_2	10	réunion				
t_3	11	congés				
t ₃ t ₄ t ₅	18	congés				
t ₅	18	réunion				

Requête conjonctive

Y a-t-il une réunion pendant mes congés ? ∃dtt' Calendrier(t, d, "congés") ^ Calendrier(t', d, "réunion")

→ Problème d'évaluation de requêtes sur données **probabilistes** :

Quelle est la **probabilité totale** de chaque réponse quand les faits sont présents ou absents **indépendamment** avec la **probabilité** indiquée?

Approche existante (intensionnelle):

	Calendrier					
	id	jour	type			
	t ₁ t ₂ t ₃ t ₄ t ₅	9 10 11 18 18	congés réunion congés congés réunion		PTIME	Formule de provenance
∃d	iunion alendri	conjonctive pendant mes er(t, d, "cong (t', d, "réunio	és")	~	t ₄ ∧ t ₅	

→ Problème d'évaluation de requêtes sur données **probabilistes** :

Quelle est la **probabilité totale** de chaque réponse quand les faits sont présents ou absents **indépendamment** avec la **probabilité** indiquée?

Approche existante (intensionnelle):

	Calendrier						
	id	jour	type				
	t ₁	9	congés	95%			
	t_2	10	réunion	20%			
	t_3	11	congés	30%			Formule de
	t ₄	18	congés	80%			
	t ₅	18	réunion	90%		PTIME	provenance
30	ne ré ltt' <mark>C</mark> a	iunion alendri	conjoncti pendant i er(t, d, "co (t', d, "réu	۰ ۱		t ₄ ∧ t ₅	

→ Problème d'évaluation de requêtes sur données **probabilistes** :

Quelle est la **probabilité totale** de chaque réponse quand les faits sont présents ou absents **indépendamment** avec la **probabilité** indiquée?

Approche existante (intensionnelle):

	Calendrier								
	id	jour	type						
	t ₁	9	congés	95%					
	t_2	10	réunion	20%		_			
	t_3	11	congés	30%			Formule de		
	t ₄	18	congés	80%				#P-difficile	Probabilité
	t ₅	18	réunion	90%		PTIME	provenance	en général _	
∃¢	ne rë Itt' <mark>C</mark> a	éunion alendri	conjoncti pendant i ier(t, d, "co (t', d, "réu	mes con ongés")	_	->	t ₄ Λt ₅		- 72%

Données probabilistes : résultats

→ J'ai montré comment exploiter la structure des données :

Théorème [Amarilli, Bourhis, Senellart, ICALP'15] L'évaluation de requêtes MSO est faisable en temps linéaire sur des données probabilistes de largeur d'arbre bornée

→ En un sens, ce résultat ne peut pas être amélioré (dichotomie) :

Théorème [Amarilli, Bourhis, Senellart, PODS'16] L'évaluation probabiliste de certaines requêtes FO est infaisable sur n'importe quelle famille constructible de graphes de largeur d'arbre non bornée

Genéraliser les résultats de [Courcelle, 1990] à la provenance

Données de largeur bornée

Calendrier							
id	jour	type					
t ₁	9	congés					
t_2	10	réunion					
t_3	11	congés					
t ₃ t ₄ t ₅	18	congés					
t_5	18	réunion					

Requête MSO

Y a-t-il une réunion pendant mes congés ? ∃dtt' Calendrier(t, d, "congés") ∧ Calendrier(t', d, "réunion")

Genéraliser les résultats de [Courcelle, 1990] à la provenance

Données de largeur bornée

Calendrier						
id	jour	type				
t ₁	9	congés				
t_2	10	réunion				
t_3	11	congés				
t ₃ t ₄ t ₅	18	congés				
t_5	18	réunion				

Requête MSO

Y a-t-il une réunion pendant mes congés ?

Edit Calendrier(t, d, "congés")

Λ Calendrier(t', d, "réunion")

Automate d'arbres



Genéraliser les résultats de [Courcelle, 1990] à la provenance

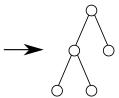
Données de largeur bornée

Calendrier jour type9 congés 10 réunion

 $egin{array}{ccccc} t_1 & 9 & {\rm cong\'es} \\ t_2 & 10 & {\rm r\'eunion} \\ t_3 & 11 & {\rm cong\'es} \\ t_4 & 18 & {\rm cong\'es} \\ t_5 & 18 & {\rm r\'eunion} \\ \end{array}$

id

Encodage en arbre



Requête MSO

Y a-t-il une réunion pendant mes congés ?

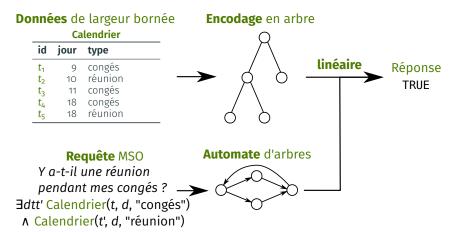
Idtt' Calendrier(t, d, "congés")

A Calendrier(t', d, "réunion")

Automate d'arbres



Genéraliser les résultats de [Courcelle, 1990] à la provenance



Genéraliser les résultats de [Courcelle, 1990] à la provenance

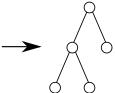
Données de largeur bornée

id

Calendrier jour type congés 10 réunion

congés 18 congés réunion

Encodage en arbre



Requête MSO

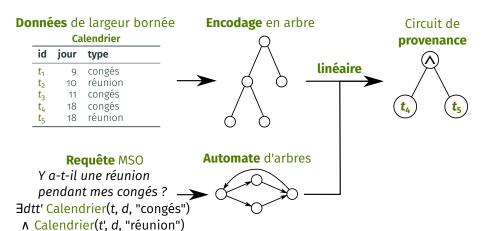
Y a-t-il une réunion pendant mes congés?

∃dtt' Calendrier(t, d, "congés") A Calendrier(t', d, "réunion")

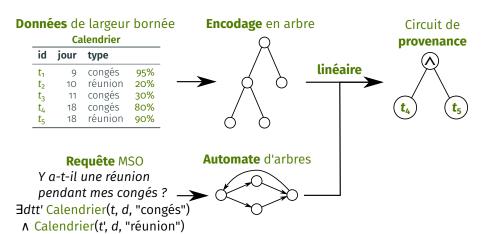
Automate d'arbres



Genéraliser les résultats de [Courcelle, 1990] à la provenance

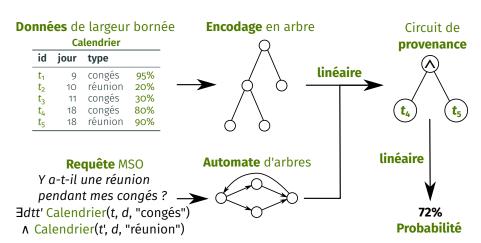


Genéraliser les résultats de [Courcelle, 1990] à la provenance



Données probabilistes : preuve de la borne supérieure

Genéraliser les résultats de [Courcelle, 1990] à la provenance



Projet de recherche

Raisonner avec la provenance sur les données du Web

• Bases de connaissances : Wikidata, YAGO, etc.

· Bases de connaissances : Wikidata, YAGO, etc.

Centre national de la recherche scientifique (Q280413)

nature de l'élément institut de recherche coordonnées géographiques + ajouter coordonnées

coordonnées géographiques 48°50'51.72"N, 2°15'50.40"E

→ 1 référence

+ ajouter

· Bases de connaissances : Wikidata, YAGO, etc.

« Centre national de la recherche scientifique »



- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.



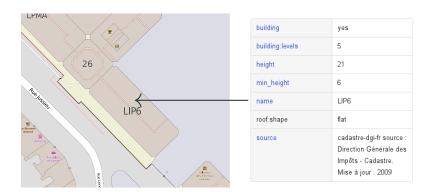
- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.



- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.
- Données géographiques : OpenStreetMaps, etc.



- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.
- Données géographiques : OpenStreetMaps, etc.



- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.
- · Données géographiques : OpenStreetMaps, etc.
- · Annotations sémantiques : Web Data Commons, etc.

Crawl Date	Winter 2014	
Total Data	64 Terabyte	
Parsed HTML URLs	2,014,175,679	
URLs with Triples	620,151,400	
Domains in Crawl	15,668,667	
Domains with Triples	2,722,425	
Typed Entities	5,516,068,263	
Triples	20,484,755,485	

- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.
- · Données géographiques : OpenStreetMaps, etc.
- · Annotations sémantiques : Web Data Commons, etc.

Objectif : combiner et utiliser ces données :

- · Répondre à des requêtes logiques complexes
- · Calculer des visualisations et des statistiques
- · Recouper des informations ou trouver des contradictions
- · Connecter ses propres données aux jeux de données existants

- · Bases de connaissances : Wikidata, YAGO, etc.
- · Données ouvertes : data.gouv.fr, etc.
- · Données géographiques : OpenStreetMaps, etc.
- · Annotations sémantiques : Web Data Commons, etc.

Objectif : combiner et utiliser ces données :

- · Répondre à des requêtes logiques complexes
- · Calculer des visualisations et des statistiques
- · Recouper des informations ou trouver des contradictions
- · Connecter ses propres données aux jeux de données existants

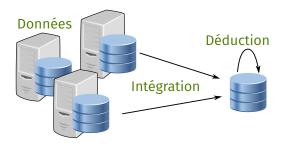
Ces données posent de nombreuses difficultés :

· hétérogènes · distribuées · incomplètes · peu fiables

Problème 1: Intégration

Données incomplètes et hétérogènes issues de sources multiples

- · Raisonnement avec des règles logiques
 - → Intégrer les différences sources
 - → Déduire les faits manquants



→ **Approches existantes** : OBDA, data integration, data exchange...

Données produites collaborativement ou par des processus faillibles

· Vandalisme

Données produites collaborativement ou par des processus faillibles

Vandalisme

Dans Wikidata

Version du 1 avril 2014 à 07:22 (remercier)

Qlsusu (discussion | contributions)

(Affirmation ajoutée : identifiant bibliothèque du Congrès (P244):

Michael Jackson tops global download sales list)

Modification suivante →

propriété / identifiant bibliothèque du Congrès

+ Michael Jackson tops global download sales list

Données produites collaborativement ou par des processus faillibles

Vandalisme

Dans Wikidata

This is your last warning. The next time you harm Wikidata, you may be blocked from editing without further notice.

Version du 1 avril 2014 à 07:22 (remercier)

Qlsusu (discussion | contributions)
(Affirmation ajoutée : identifiant bibliothèque du Congrès (P244):
 Michael Jackson tops global download sales list)
 Modification suivante →

propriété / identifiant bibliothèque du Congrès

+ Michael Jackson tops global download sales list

Données produites collaborativement ou par des processus faillibles

Vandalisme

Controverses

Dans Wikidata



Données produites collaborativement ou par des processus faillibles

Vandalisme

· Controverses · Extraction

Dans YAGO



Macquarie Group Limited MACQUARIE Public Type Founded

Sydney, New South Wales, Australia (1970)AUD A\$8.1 billion Revenue $(2014)^{[1]}$

Données produites collaborativement ou par des processus faillibles

Vandalisme

Controverses

Extraction

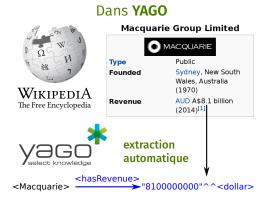


Données produites collaborativement ou par des processus faillibles

Vandalisme

Controverses

Extraction



→ Approches existantes : truth finding, data cleaning, data repair...

Objectif

L'intégration et la fiabilité vont ensemble pour les données du Web :

- · Les données d'une source peuvent provenir d'autres sources
- · Il faut étudier la fiabilité des résultats du raisonnement
- · Les règles d'intégration elles-mêmes ne sont pas fiables

Objectif : Intégrer des données et raisonner sur ces données en conservant des informations pour estimer leur fiabilité

Approche : Définir et utiliser la provenance des données

• 1. Exploiter la provenance • 2. Propager la provenance

Exploiter la provenance des données

Wikidata: Plus de 40M faits ont une source (près de 50%)





ShonagonBot (discussion | contributions)
(Ajout d'une référence à une affirmation : Identifiant BnF
(1268): 118629315)
Modification suivante →

propriété / identifiant BnF : 118629315 / référence

importé de : data.bnf.fr
date de consultation :
26 août 2015
Horodatage +2015-08-26700.00.00Z
Zone horaire +00.00
Calendrier Grégorien
affirmé dans : data.bnf.fr

Exploiter la provenance des données

Wikidata : Plus de 40M faits ont une source (près de 50%)



Version du 30 août 2015 à 12:28 (restaurer) (annuler) ShonagonBot (discussion | contributions) (Ajout d'une référence à une affirmation : identifiant BnF (#288): 118629315) Modification suivante → propriété / identifiant BnF : 118629315 / référence



OpenStreetMaps: >40M points et >120M voies avec source (>35%)





source cadastre-dgi-fr source :
Direction Générale des
Impôts - Cadastre.
Mise à jour : 2009

Exploiter la provenance des données

Wikidata : Plus de 40M faits ont une source (près de 50%)



Version du 30 août 2015 à 12:28 (restaurer) (annuler) ShonagonBot (discussion | contributions) (Ajout d'une référence à une affirmation : identifiant BnF (12:28): 118629315) Modification suivante → propriété / identifiant BnF : 118629315 / référence



OpenStreetMaps: >40M points et >120M voies avec source (>35%)





Way: LIP6 (42741440)

cadastre-dgi-fr source :
Direction Générale des
Impôts - Cadastre.
Mise à jour : 2009

Autres indices : Historique des faits, utilisateurs qui ont édité

Propager la provenance des données

→ Comment définir la **provenance** d'une réponse **certaine**?

Données:

id	jour	type
t ₁	9	congés
t ₂	10	réunion
t ₃	11	congés

Requête logique :

Y a-t-il des réunions pendant mes congés?

Règle logique :

Je ne reviens pas pour **un seul jour**

Provenance:

 $t_1 \wedge t_2 \wedge t_3 \wedge r\grave{e}gle$?

Propager la provenance des données

→ Comment définir la **provenance** d'une réponse **certaine**?

Données:		s :	Requête logique :	Règle logique :	Provenance:
id	jour	type	Y a-t-il des réunions	Je ne reviens pas	$t_1 \wedge t_2 \wedge t_3 \wedge règle$
	9	congés réunion	pendant mes congés?	pour un seul jour	11//12//13//regie
t ₂			peridunt mes conges :	pour un seut jour	
t ₂	11	conges			

→ Généraliser les **semianneaux de provenance** au raisonnement?

Propager la provenance des données

→ Comment définir la **provenance** d'une réponse **certaine**?

Données:		s:	Requête logique :	Règle logique :	Provenance:
id	jour	type	Y a-t-il des réunions	Je ne reviens pas	$t_1 \wedge t_2 \wedge t_3 \wedge règle$
t ₁ t ₂ t ₃	10	congés réunion congés	pendant mes congés?	pour un seul jour	11//12//13//regie

- → Généraliser les **semianneaux de provenance** au raisonnement?
- → Comment calculer **efficacement** cette provenance et la représenter de façon **concise**, selon le langage de règles?

Programme de recherche

Raisonner avec la provenance sur les données du Web pour l'intégration et la fiabilité

- 1. **Fondations**: Provenance symbolique pour le raisonnement
 - · Définition abstraite à différents niveaux d'expressivité
 - · Calcul et représentation efficace
- 2. Propager des relations de fiabilité à travers la provenance
- 3. Calculer des confiances quantitatives et probabilistes
- 4. **Long terme :** Réviser les jugements sur les sources primaires avec des retours utilisateurs, en remontant la provenance

Intégration

Intégration

Équipe LINKS de CRIStAL, Lille

Thématiques: Web des données et intégration

Expertise: Automates (S. Tison), agrégation (J. Niehren),

Collaboration: Avec P. Bourhis

Équipe Automates et applications de l'IRIF, Paris

Thématiques: Informatique fondamentale (large champ),

théorie des bases de données

Expertise: Données incertaines (C. Sirangelo et A. Gheerbrant)

Équipe GraphIK du LIRMM, Montpellier

Thématiques: Raisonnement sous contraintes expressives

Expertise: Règles existentielles (J.-F. Baget, M.-L. Mugnier) et

logiques de description (M. Bienvenu)

Résumé

Projet : Raisonner avec la provenance sur les données du Web pour l'intégration et la fiabilité

- 1. **Fondations**: Provenance symbolique pour le raisonnement
 - · Définition abstraite à différents niveaux d'expressivité
 - · Calcul et représentation efficace
- 2. Propager des relations de fiabilité à travers la provenance
- 3. Calculer des confiances quantitatives et probabilistes
- 4. **Long terme :** Réviser les jugements sur les sources primaires avec des retours utilisateurs, en remontant la provenance

Thèse: Tirer parti de la structure des données incertaines [ICALP'15], [LICS'15], [PODS'16]; après-thèse [IJCAl'16]

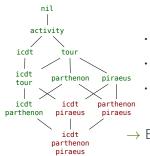
Pré-doc: À Tel Aviv [ICDT'14] et à Oxford [IJCAI'15]

Crowdsourcing

Fouille de données : Trouver des motifs fréquents

Crowdsourcing : Poser des questions à la foule

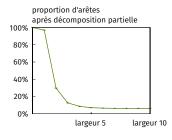
→ Questions à la foule : quels **ensembles d'objets** sont **fréquents?**



- Utiliser une taxonomie sur les objets
- · Les ensembles forment un treillis distributif
- Compromis entre le coût des questions posées et le coût de calculer quelles questions poser
- → Bornes de **complexité** sur ce problème

Applicabilité pratique

 Travail avec S. Maniu : les jeux de données réels peuvent être partiellement décomposés en arbre



Décomposition **partielle** en arbre du **graphe OSM de Paris**

- · 4.3 M nœuds et 5.4 M arêtes
- Largeur totale \leq 521
- Stage de M. Monet : les méthodes à base d'automates peuvent être implantées en pratique
- Travail avec M. Monet : compilation efficace en automates pour des langages restreints de requêtes

Références



COURCELLE, Bruno (1990). "The Monadic Second-Order Logic of Graphs. I. Recognizable Sets of Finite Graphs". In: Inf. Comput.