

5. Gestion des données et vie privée du citoyen urbain

- Architectures de gestion de données face au respect de la vie privée
- Gestion de la vie privée dans les réseaux sociaux mobiles
- Privacy-by-design: gestion de données confinées (puces et capteurs)
- Gestion de la vie privée dans les applications mobiles participatives
- **Traitements de données globaux respectueux de la vie privée des participants**

Nicolas Anciaux

VILLES INTELLIGENTES : DÉFIS TECHNOLOGIQUES ET SOCIÉTAUX

Cette dernière séquence présente les techniques de calcul distribué respectueuses de la vie privée et qui peuvent servir à calculer des résultats obtenus à partir des données personnelles de grand nombre d'individus.

Calculs globaux sécurisés à large échelle

- **Traitement à large échelle**

- Calcul statistiques, fouille de données, Big Data, etc.

- **Hypothèses**

- Données de base détenues par les individus
- Seul le résultat du calcul doit être révélé
- L'infrastructure de support (calcul, comm.) n'est pas de confiance

- **Participants et l'infrastructure**

- Honnêtes, semi honnêtes, malveillants

2

Quels sont les traitements globaux que l'on souhaite réaliser ?

Du **traitement à grande échelle**, comme des calculs statistiques par exemple pour réaliser des études épidémiologiques, de la fouille de données pour classifier, calculer des profils ou tout type de traitements Big Data basés sur les techniques map reduce et bien d'autres calculs encore.

Les **hypothèses** que l'on fait dans cette séquence sont :

- les données de base du calcul sont détenues par les individus
- seul le résultat du calcul doit être révélé.
- L'infrastructure extérieure qui peut servir au calcul ou qui sert aux échanges de données entre les participants ne doit pas avoir accès aux données de base qui permettent de calculer le résultat.

Les **participants et l'infrastructure** peuvent donc être considérés comme étant :

- honnêtes, c'est-à-dire qu'ils suivent le protocole,
- semi-honnête, ils suivent aussi le protocole mais cette fois, ils tentent d'inférer de l'information
- malveillants et dans ce cas-là, ils ne suivent pas le protocole et essayent essentiellement d'inférer de l'information.

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- **Bases de données chiffrées centralisées (DBaaS)**
 - L'hébergeur n'a pas les clés mais le SGBD peut calculer sur le chiffré
 - Comment calculer sur le chiffré ?

3

La première solution qu'on peut imaginer, c'est d'utiliser du chiffrement.

Chaque participant va chiffrer ses données et va les transmettre à destination d'une base de données centralisée qui va effectuer le calcul sur le chiffré.

La question qui se pose, c'est comment le serveur peut-il calculer le résultat si les données sont chiffrées et que bien sûr, le serveur ne dispose pas de la clé qui permet de les déchiffrer?

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- Bases de données chiffrées centralisées (DBaaS)
- Chiffrement homomorphe
 - Préserve certaines opérations. Ex: RSA : $C(m1) \times C(m2) = C(m1 \times m2)$
 - Comment réaliser n'importe quel calcul ?

4

On peut utiliser du chiffrement homomorphe. C'est un mode de chiffrement qui **préserve certaines opérations** qui peuvent être réalisées directement sur les données chiffrées. Par exemple si on chiffre des données avec un algorithme RSA, on va préserver la multiplication de manière à ce que si on multiplie entre eux 2 nombres chiffrés, le résultat donne un nombre chiffré qui une fois déchiffré est le résultat de la multiplication sur les nombres en clair.

Mais peut-on réaliser de cette manière n'importe quel calcul ? Pas encore, du moins pas à l'heure actuelle parce qu'il faudrait du **chiffrement complètement homomorphe** qui préserve suffisamment d'opérations pour pouvoir implémenter n'importe quel traitement.

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- Bases de données chiffrées centralisées (DBaaS)
- Chiffrement homomorphe
- Chiffrement complètement homomorphe [Gent09]
 - Évalue toute fonction. Mais coût en performances incroyablement élevé

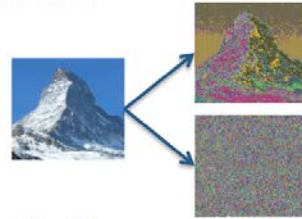
5

Il y a des résultats théoriques là-dessus, mais qui ne sont pas encore exploitables en pratique, car incroyablement coûteux, et en tout cas pas exploitables dans le cadre de la gestion de données.

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- Bases de données chiffrées centralisées (DBaaS)
- Chiffrement homomorphe
- Chiffrement complètement homomorphe [Gent09]
- Utiliser du chiffrement déterministe
 - Préserve l'égalité, l'ordre, etc. Ex. $m_1 = m_2 \Leftrightarrow C(m_1) = C(m_2)$
 - Inférences ?



6

Pour effectuer les opérations utiles au traitement base de données, il faudra donc utiliser des modes de chiffrement laissant la possibilité de réaliser des traitements de données quitte à dégrader la sécurité. Par exemple, un mode de chiffrement qui est largement utilisé pour préserver les traitements en bases de données est le chiffrement déterministe.

L'effet en est montré sur cette image :

- dans celle du haut chaque pixel est chiffré en déterministe,
- dans celle du bas en non déterministe.

L'avantage avec l'image du haut est que les pixels restent indexables dans la base de données. Ainsi rechercher du bleu dans cette image revient à rechercher tous les pixels ayant pour valeur la valeur chiffrée de bleu.

Avec ce type de chiffrement, l'indexation est permise, mais la distribution de données est bien sûr visible.

Si on connaissait cette distribution, il est facile de renverser le chiffrement.

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- Bases de données chiffrées centralisées (DBaaS)
- Chiffrement homomorphe
- Chiffrement complètement homomorphe [Gent09]
- Utiliser du chiffrement déterministe
- Minimiser les inférences
 - Restreindre les propriétés vs. l'usage [AdaPopa'11]
 - Etiquettes indistinguables aux données [Hacigums'02]

7

Le problème est donc aussi de chiffrer mais en se protégeant des inférences. Certaines techniques permettent par exemple de **favoriser du chiffrement opaque** à chaque fois que c'est possible et d'autres permettent d'attacher aux données des **étiquettes qui sont statistiquement indistinguables** et qui permettent quand même au serveur d'évaluer certains calculs.

Publier des jeux de données chiffrés

... et réaliser le calcul directement sur les données chiffrées

- Bases de données chiffrées centralisées (DBaaS)
- Chiffrement homomorphe
- Chiffrement complètement homomorphe [Gent09]
- Utiliser du chiffrement déterministe
- Minimiser les inférences

Limites actuelles :

- Balance performance/sécurité
- Inférences
- Utilisateurs avec des droits différenciés

8

Ces techniques ont actuellement de nombreuses limites :

- le **compromis** fait sur la **sécurité** pour conserver de bonnes **performances** est fort, car les **inférences** sont possibles.
- Si les utilisateurs ont des droits différenciés, il faudra chiffrer les données avec des clés de chiffrement différentes, ce qui rend la gestion très complexe.

Publier des jeux de données anonymes

... et réaliser le calcul directement sur les données anonymes

- **Anonymisation de données sensibles**

- Dégrader les informations pour empêcher de retrouver celles d'un individu donné
- Pseudonymat → k-anonymat [Sweeney02] → l-diversité [Kifer07] → t-fermeture [Li07] ...
- Garanties tangibles mais pas « formelle »

9

Une deuxième alternative consiste à anonymiser les données de chacun et à réaliser le calcul sur les données anonymes.

Une première solution pour obtenir l'anonymat est de généraliser ou supprimer certaines informations jusqu'à éviter de pouvoir identifier celles qui correspondent à un individu donné.

Mais plus le jeu de données devient anonyme et moins il est utilisable.

Un simple pseudonymat où les identifiants présents dans le jeu données sont remplacés par un pseudonyme ne suffit pas.

Le **k-anonymat** va plus loin et propose de généraliser aussi les éléments de jeu données qui sont perçus comme des quasi identifiants et qui pris en combinaison, pourraient identifier une personne en particulier.

La **l-diversité** propose de ne jamais associer une seule et même valeur sensible pour un même jeu de quasi identifiants.

La **t-fermeture** va encore plus loin en cherchant à respecter les distributions connues pour les données sensibles.

Publier des jeux de données anonymes

... et réaliser le calcul directement sur les données anonymes

- **Anonymisation de données sensibles**
- **Garantie différentielle d'anonymat** [Dwork06]
 - Ajouter du bruit et faire un calcul approché
 - Que les données d'un individu soient là ou pas, le résultat est (presque) le même
 - Comment s'abstraire d'un tiers de confiance qui constitue le jeu anonyme ?

10

Une seconde technique consiste à introduire de fausses données de manière à masquer le fait qu'un individu appartient ou n'appartient pas au jeu de données.

Ainsi, un même calcul réalisé sur 2 jeux de données qui ne diffèrent que par les données d'un seul individu donnera le même résultat.

Les garanties offertes sont ainsi plus formelles qu'avec les techniques précédentes.

Mais un autre problème se pose: comment générer un tel jeu de données anonymes?

Publier des jeux de données anonymes

... et réaliser le calcul directement sur les données anonymes

- **Anonymisation de données sensibles**
- **Garantie différentielle d'anonymat** [Dwork06]
- **Produire un jeu anonyme sans centraliser les données ?**
 - Problème : usage (anonymise-puis-intègre) << usage (intègre-puis-anonymise)
 - ➔ Dans ce cas précis, possible avec des techniques cryptographiques [Mohammed10]

11

Si chaque individu anonymise lui-même ses données localement, l'usage qui pourra être fait des données sera très dégradé.

Pour anonymiser correctement, il faut confronter les données de l'ensemble de tous les individus pour produire le jeu anonyme

Publier des jeux de données anonymes

... et réaliser le calcul directement sur les données anonymes

- Anonymisation de données sensibles
- Garantie différentielle d'anonymat [Dwork06]
- Produire un jeu anonyme sans centraliser les données ?
 - Problème : usage (anonymise-puis-intègre) << usage (intègre-puis-anonymise)
 - ➔ Dans ce cas précis, possible avec des techniques cryptographiques [Mohammed10]

Plus généralement : peut-on produire le résultat d'un calcul sans en exposer les données de base ?

12

. Ca revient donc au problème présenté précédemment, celui de calculer un résultat sans révéler les données de base qui permettent le calcul.

Calculer sans centraliser : oui !

- **Calcul multipartite sécurisé (CMS) : calculer sans révéler les entrées**
 - Exemple : problème du millionnaire [Yao82]
 - Coût d'un calcul générique exponentiel dans le nombre d'entrées
- **Fouille de données distribuée respectueuse de la vie privée**
 - Kit de techniques issues du CMS et adaptées à la fouille de données [Clifton02]
 - ✓ Somme, union, nombre d'éléments d'une intersection, produit scalaire
 - ✓ Calculs rendus possibles : règles d'association, classification

13

C'est un problème qui est étudié en cryptographie dans le cadre du **calcul multipartite sécurisé** (CMS). Mais le coût d'un calcul générique est exponentiel dans le nombre d'entrées.

Certaines **techniques de fouille de données** peuvent toutefois être dérivées de ces techniques de calcul multipartite sécurisé et **permettent de calculer des règles d'associations ou de faire de la classification de manière sécurisée**.

Calculer sans centraliser : oui !

- **Calcul multipartite sécurisé (CMS) : calculer sans révéler les entrées**
 - Exemple : problème du millionnaire [Yao82]
 - Coût d'un calcul générique exponentiel dans le nombre d'entrées
- **Fouille de données distribuée respectueuse de la vie privée**
 - Kit de techniques issues du CMS et adaptées à la fouille de données [Clifton02]
 - ✓ Somme, union, nombre d'éléments d'une intersection, produit scalaire
 - ✓ Calculs rendus possibles : règles d'association, classification

**Limites :
performance / généricité**

14

Mais à l'heure actuelle, ces techniques cryptographiques ne permettent pas de répondre à la plupart des problèmes de gestion de données.

Nouvel angle d'attaque : matériel sécurisé

- **Introduction de matériel sécurisé (de confiance)**
 - Simplifie drastiquement les traitements CMS (et les traitements BD associés)
 - Applications : CMS avec matériel sécurisé [Katz07, GIS+10], calculs base de données [Allard10], [To14]

15

Un nouvel angle d'attaque émerge actuellement grâce à l'introduction de matériels sécurisés ou de matériels de confiance qui simplifient largement la problématique classique du calcul multipartite sécurisé.

Nouvel angle d'attaque : matériel sécurisé

- Introduction de matériel sécurisé (de confiance)
- **Matériel sécurisé côté serveur** : utilisé comme tiers de confiance centralisé
 - Co-processeur résistant aux attaques physiques
 - ✓ Ex: IBM 4764 PCI-X Cryptographic Coprocessor [IBM]
 - Un SGBD complet peut y être embarqué [Bajaj14]

16

Par exemple un **co-processeur sécurisé** qui ne peut pas être observé peut être placé côté serveur pour collecter des données chiffrées - donc celles des participants - et les déchiffrer en interne pour calculer un résultat qui seul sera exposé. La solution Trusted DB a été proposée dans cette optique.

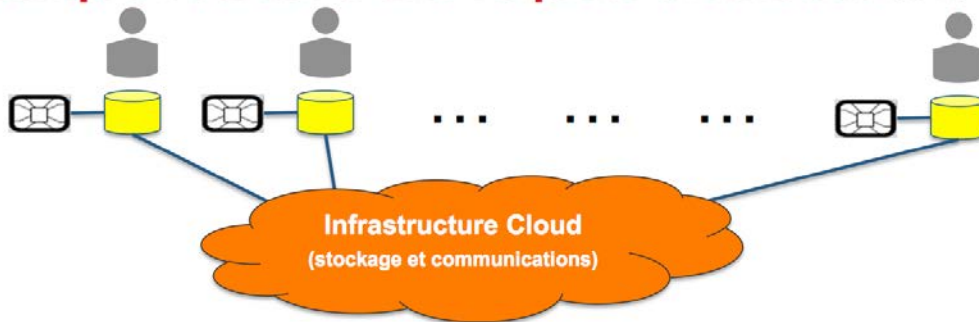
Nouvel angle d'attaque : matériel sécurisé

- **Introduction de matériel sécurisé (de confiance)**
- **Matériel sécurisé côté serveur** : utilisé comme tiers de confiance centralisé
- **Matériel sécurisé côté client** : nœud de calcul (de confiance, très contraint)
 - Carte SIM, token ou carte sécurisée, ...
 - Des traitements de données simples sont embarqués [Anciaux14]
.... en support à des calculs distribués [Katz07, Javinen10, Allard10]
 - ✓ Ex : Map-Reduce et agrégats SQL [To14]

17

` Autre exemple, du matériel sécurisé peut être présent côté client, il peut s'agir d'une carte SIM embarquée dans un téléphone portable ou de tout autre objet sécurisé à disposition de l'individu. Comme nous l'avons vu lors de la troisième séquence de cette semaine, les traitements peuvent être embarqués dans ce type de dispositif. Si les dispositifs peuvent s'échanger entre eux des données via une infrastructure, le tout peut être vu comme une plate-forme distribuée et sécurisée. À l'heure actuelle, des calculs d'agrégats SQL et des traitements Map reduce, par exemple, peuvent être réalisés de cette manière.

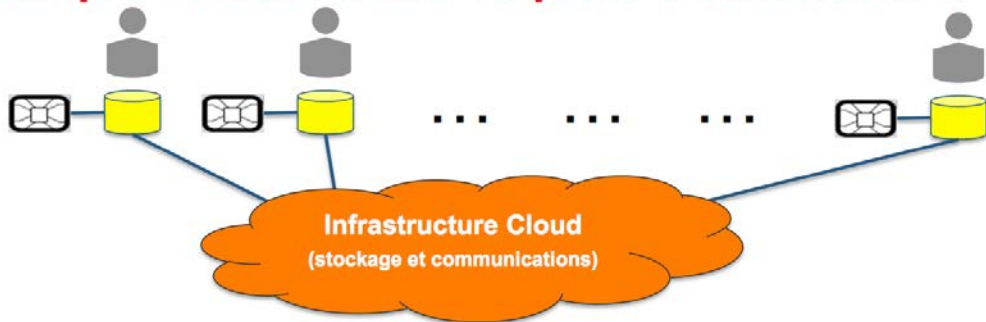
Exemple : calcul d'une requête distribuée [To14]



18

De façon très synthétique, voilà comment fonctionne ce type de calcul.
Supposons un grand nombre de citoyens disposant chacun de leurs données et d'un dispositif personnel sécurisé. Ces citoyens doivent être capables d'échanger des données au travers d'une infrastructure.

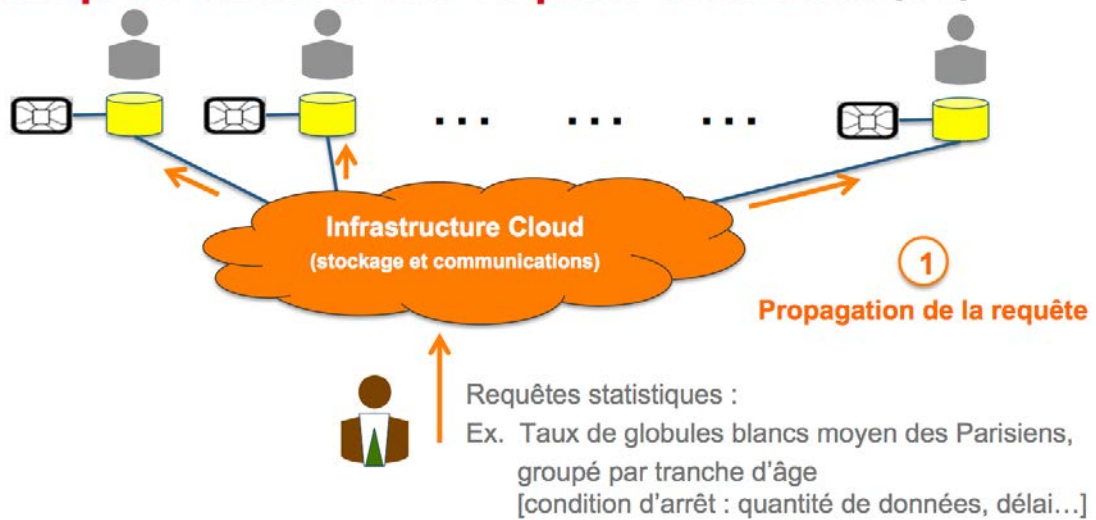
Exemple : calcul d'une requête distribuée [To14]



Requêtes statistiques :

Taux de globules blancs moyen des Parisiens,
groupé par tranche d'âge
[condition d'arrêt : quantité de données, délai...]

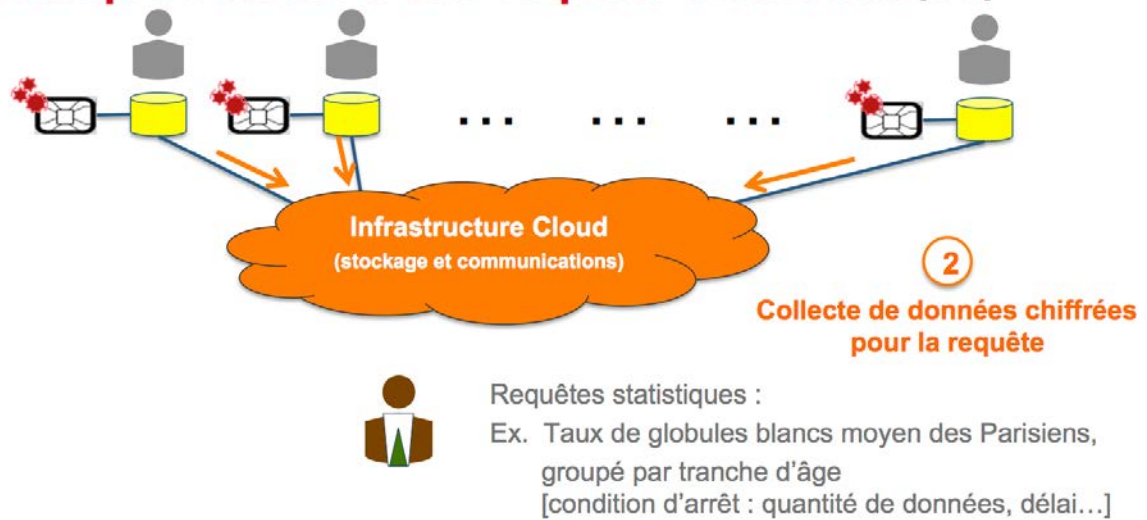
Exemple : calcul d'une requête distribuée [To14]



20

Dans une première phase, la requête que cet individu va poser est chiffrée et transmise aux participants.

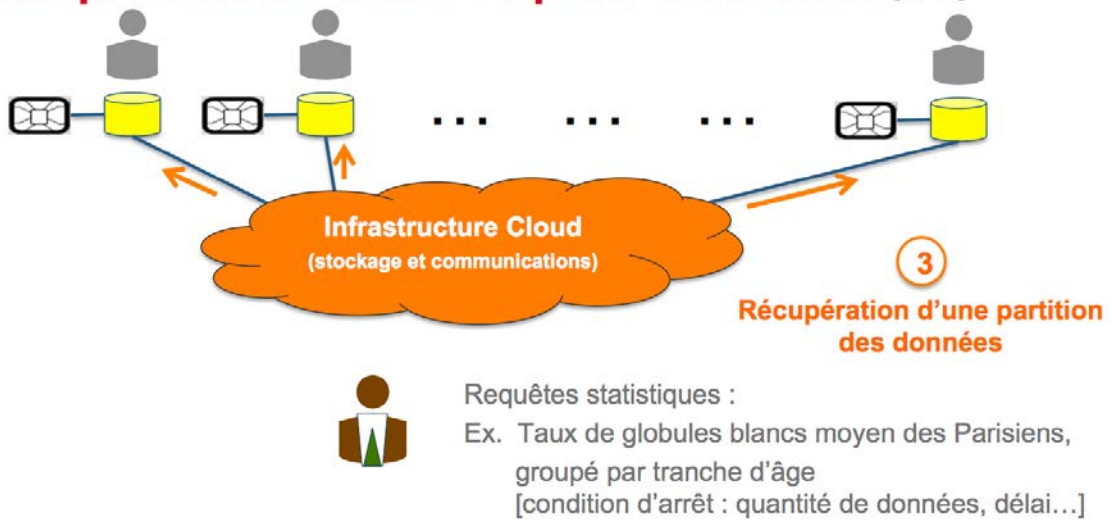
Exemple : calcul d'une requête distribuée [To14]



21

Le matériel sécurisé déchiffre la requête, évalue localement l'ensemble des données autorisées et pertinentes pour cette requête et transmet ces données chiffrées à l'infrastructure.

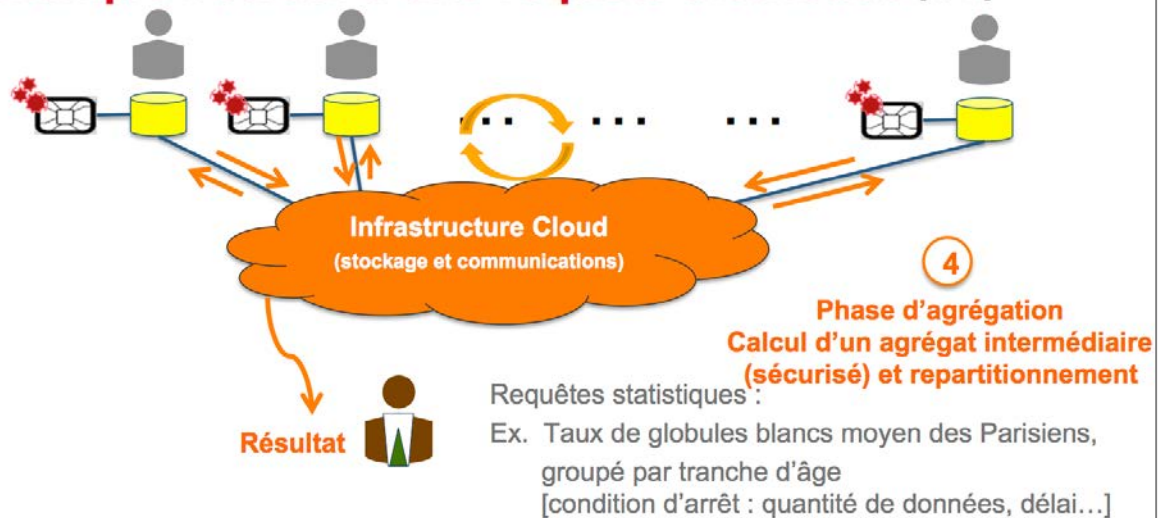
Exemple : calcul d'une requête distribuée [To14]



22

Les données sont partitionnées par l'infrastructure et chaque participant en récupère une partition.

Exemple : calcul d'une requête distribuée [To14]



Une phase de calcul peut commencer. Chaque dispositif sécurisé réalise un calcul partiel sur sa partition et renvoie les résultats chiffrés à l'infrastructure.

Un mécanisme itératif s'engage jusqu'à ce que le résultat final puisse être construit et le résultat final une fois construit peut être récupéré et déchiffré par l'individu ayant posé la question.

Conclusion

- **Techniques de calcul distribué respectueux de la vie privée**
 - Anonymat, garanties différentielles, cryptographie, matériel sécurisé
 - Tendent vers des solutions plus génériques, moins coûteuses et plus protectrices
- **Bases de données chiffrées**
- **Calcul multipartite sécurisé et application à la fouille de données**
- **Introduction de matériel sécurisé qui simplifie le problème**
- **Impact en pratique reste limité**
 - (?) Lacune d'origine : manque de généricité, coût (performances) important
 - (?) Valeur intrinsèque des données, modèles d'affaires

24

En conclusion, de nombreuses techniques existent actuellement pour calculer des résultats globaux dans le respect de la vie privée.

Ces techniques sont basées sur de l'**anonymat**, du **chiffrement**, du **calcul multipartite sécurisé** ou du **matériel sécurisé**.

L'impact de ces techniques en pratique reste encore limité :

- (?) Lacune d'origine : manque de généricité, coût (performances) important
- (?) Valeur intrinsèque des données, modèles d'affaires

Cependant, l'effort normatif actuel, notamment conduit en Europe avec la refonte de la directive de 95, et les effets de bord de l'affaire Snowden nous permettent d'espérer une accélération de l'utilisation de ce genre de techniques au bénéfice des villes intelligentes.

Références 1/2

- [Sweeney02] Sweeney, L. (2002). Achieving k-anonymity privacy protection using generalization and suppression. *International Journal of Uncertainty, Fuzziness and Knowledge-Based Systems*, 10(5).
- [Kifer07] Machanavajjhala, A., Kifer, D., Gehrke, J., & Venkatasubramanian, M. (2007). l-diversity: Privacy beyond k-anonymity. *ACM Transactions on Knowledge Discovery from Data (TKDD)*, 1(3).
- [Li07] Li, N., Li, T., & Venkatasubramanian, S. (2007, April). t-closeness: Privacy beyond k-anonymity and l-diversity. In *ICDE 2007*.
- [Yao82] Yao, A.C.: Protocols for secure computations. In *Annual Symposium on Foundations of Computer Science, FOCS*, 1982.
- [Dwork06] Dwork, C. (2011). Differential privacy. In *Encyclopedia of Cryptography and Security*.
- [Bambauer14] Bambauer, J., Muralidhar, K., & Sarathy, R. (2014). Fool's Gold: an Illustrated Critique of Differential Privacy. *Vand. J. Ent. & Tech. L.*, 16, 701.
- [Clifton02] Clifton, C., Kantarcioglu, M., Vaidya, J., Lin, X., & Zhu, M. Y. (2002). Tools for privacy preserving distributed data mining. *ACM Sigkdd Explorations Newsletter*, 4(2).
- [To14] To, Q. C., Nguyen, B., & Pucheral, P. (2014). Privacy-Preserving Query Execution using a Decentralized Architecture and Tamper Resistant Hardware. In *EDBT'14*.

Références 2/2

- [Hacigums'02] Hacigümüş, H., Iyer, B., Li, C., & Mehrotra, S. (2002). Executing SQL over encrypted data in the database-service-provider model. *ACM SIGMOD*.
- [AdaPopa'11] Ada Popa, R., Redfield, C. M. S., Zeldovich, N. and Balakrishnan, H. (2011). CryptDB: Protecting Confidentiality with Encrypted Query Processing. *ACM SOSP*.
- [Mohammed10] Mohammed, N., Fung, B., Hung, P. C., & Lee, C. K. (2010). Centralized and distributed anonymization for high-dimensional healthcare data. *ACM TKDD*, 4(4).
- [IBM] IBM 4764 Model 001 PCI-X Cryptographic Coprocessor. Lien : https://www-03.ibm.com/security/cryptocards/pdfs/4764-001_PCIX_Data_Sheet.pdf
- [Bajaj14] Bajaj, S., & Sion, R. (2014). TrustedDB: A trusted hardware-based database with privacy and data confidentiality. In *IEEE TKDE*, 26(3).
- [Javinen10] Jarvinen, K., Kolesnikov, V., Sadeghi A-R., Schneider, T.: Embedded SFE: Offloading Server and Net-work Using Hardware Tokens. In *Financial Cryptography and Data Security* (2010)
- [Katz07] Katz, J.: Universally Composable Multi-party Computation Using Tamper-Proof Hardware. In *EUROCRYPT*, 2007.
- [Allard10] Allard, T., Anciaux, N., Bouganim, L., Guo, Y., Le Folgoc, L., Nguyen, B., Pucheral, P. & Yin, S. (2010). Secure personal data servers: a vision paper. In *PVLDB*, 3(1-2).