Sémantique et extension du langage C2QL pour la composition de techniques protégeant la confidentialité dans le nuage

Santiago Bautista

Juillet 2017

Résumé

Des applications de tout genre manipulent des personnelles de ses utilisateurs et utilisent le cloud pour s'exécuter ou s'héberger. Différentes techniques existent pour protéger la confidentialité de ces données. Pendant ce stage on a étudié la sémantique et prouvé les propriétés algébriques d'un langage permettant de décrire efficacement la composition de plusieurs de ces techniques, comme la fragmentation et le chiffrement.

Mots clés: privacy, cloud-computing, semantics, proof of correctness, algebraic laws, fragmentation, optimisation

Table des matières

1	Intr	oduction	2
2 Contexte		itexte	2
	2.1	De l'importance de composer les techniques de protection	2
	2.2	Un langage pour décrire la composition : C2QL	4
	2.3	Utiliser C2QL : commuter les opérateurs	5
3 Contribution		6	
	3.1	Établir des définitions	6
	3.2	Compléter les propriétés	7
	3.3	Prouver les propriétés	7
	3.4	Optimiser les requêtes	7
4	Travail futur		7
5	5 Conclusion		7

1 Introduction

De plus en plus de logiciels sont développés pour être exécutés dans le cloud, et ses logiciels, de quelque sorte qu'ils soient (messagerie, gestion d'agenda personnel, commande de pizza ou reconnaissance vocale) traitent des données personnelles, qu'ils doivent donc protéger.

Une des propriétés qui doit être garantie dans la protection des données personnelles est la confidentialité. Différentes techniques existent pour protéger la confidentialité des données, comme par exemple le chiffrement et la fragmentation.

Dans sa thèse de 2016, Ronan Cherrueau montre que chacune de ces techniques a des avantages et des inconvénients, mais qu'en composant les différentes techniques ensemble on peut profiter de tous les avantages de ces techniques en éliminant la plupart des inconvénients. Il développe donc un langage, nommé C2QL (pour *Cryptographic Compositions for Query Language*), qui permet de décrire une telle composition de techniques de sécurisation des données pour en vérifier la correction et raisonner plus facilement.

Le langage se présente comme un ensemble de fonctions que l'on peut composer entre elles. Parmi ces fonctions, il y a les fonctions classiques pour faire des requêtes dans des bases de données, telles que la projection et la sélection, tout comme des fonctions décrivant la protection des données, comme le chiffrement ou la fragmentation.

Un des intérêts de ce langage est que, pour décider comment protéger les données des utilisateurs, le développeur peut suivre un processus simple en trois étapes. D'abord, le développeur écrit les requêtes *en ne tenant compte* ni du fait que le programme s'exécute dans le nuage, ni des mécanismes pour le protéger. Puis, il compose sa requête avec les fonctions de protection nécessaires (qui dépendent du problème en particulier, des contraintes de confidentialité spécifiques) pour avoir une requête sécurisée. Finalement, le développeur utilise des lois de commutation entre les différentes fonctions pour optimiser sa requête sécurisée.

Par conséquent, disposer de lois qui indiquent à quelles conditions les différentes fonctions du langage commutent est très important.

Or, si dans sa thèse R. Cherrueau donne la plupart de ces lois, il n'a pas eu le temps de les démontrer, ni de contempler tous les cas de figure.

C'est pourquoi, pendant ce stage, j'ai complété l'ensemble de lois fournies (section 3.2) dans la thèse de Ronan et j'ai formalisé la sémantique des différentes fonctions (section 3.1) pour ensuite démontrer la correction de ces lois (section 3.3).

Une description du contexte dans lequel s'inscrit ce stage est donnée à la section 2, et une discussion sur les aspects qui n'ont pas été traités est donnée à la section 4.

2 Contexte

2.1 De l'importance de composer les techniques de protection

De plus en plus d'applications, en particulier les applications web et les applications pour téléphone portable, cherchent à utiliser le nuage, soit pour stocker du code ou des données, soit pour faire des calculs, voir les deux à la fois.

En effet, le nuage peut offrir des services (que ce soit sous forme d'infrastructure, de plateforme ou de logiciel) disposant d'une forte disponibilité et faciles à redimensionner.

Autrement dit, pouvoir utiliser le nuage est devenu un enjeu de la conception logicielle, à cause des avantages de disponibilité et redimensionnement que cela offre.

Mais ce n'est pas le seul enjeux de la conception logicielle.

Vu que la plupart de ces applications manipulent des données personnelles, garantir la *confidentialité* des données est également un enjeux de ces applications là ; tout comme le sont les *performances* pour garantir une meilleure expérience à l'utilisateur de l'application.

Dans sa thèse, R. Cherrueau s'intéresse à trois techniques particulières utilisées dans le développement logiciel et regarde comment elles interagissent avec les trois enjeux cités plus haut. Ces trois techniques, qu'on va décrire brièvement, sont le *chiffrement*, la *fragmentation verticale* et l'exécution de l'application *chez l'utilisateur*.

Le chiffrement Lorsqu'il est bien utilisé, le chiffrement permet de garantir la confidentialité des données de l'utilisateur. De plus, dans certains cas, des calculs peuvent être faits sur les données chiffrées. On appelle chiffrement homomorphe un chiffrement avec lequel on peut effectuer des calculs avec les données chiffrées. Les chiffrement homomorphes totaux, comme celui de Gentry (référence à ajouter) sont pour l'instant trop contraignants pour pouvoir être utilisés dans la plupart des applications, mais les chiffrements homomorphes partiels, c'est à dire les chiffrement avec lesquels on peut effectuer certaines opérations sur les données chiffrées, peuvent se révéler très utiles. C'est le cas des chiffrements déterministes (dont les chiffrements symétriques) qui sont des chiffrement homomorphes partiels, permettant le test d'égalité.

Dans tous les cas, le chiffrement implique un surcoût en terme de calculs, donc diminue les performances. Dans certains cas, il améliore la confidentialité et permet l'utilisation du nuage.

L'exécution côté client Si le programme était exécuté entièrement par la machine de l'utilisateur, cela serait à la fois bon pour la confidentialité (car les données ne seraient pas du tout exposées aux risques liés à l'utilisation du nuage et du réseau) et pour les performances, puisque, à moins de traiter une trop grande quantité de données dans un calcul hautement parallélisable, les performances du cloud sont moins bonnes que celles des machines des utilisateurs. Par contre, l'exécution côté client ne permet pas de profiter des avantages du cloud.

La fragmentation verticale consiste à séparer les différentes données que manipule le programme entre deux clouds n'ayant aucun rapport entre eux (à deux endroits géographiques différents, gérés par des entités différentes, etc...). Ceci permet de protéger celles des données personnelles qui sont constituées d'une association de deux données. Par exemple, dans une application stockant un ensemble de rendez-vous, l'association (date, lieu) doit être protégée, car une personne malveillante ayant accès à ces deux informations là à la fois pourrait suivre l'utilisateur de l'application.

La fragmentation verticale contribue à protéger la confidentialité, mais d'une façon souvent moins forte que celles du chiffrement ou de l'exécution côté client. Par contre, chacun des fragments de données qu'elles génère peuvent être opérés séparément, en introduisant ainsi, lorsque le programme le permet, une dose de parallélisation supplémentaire qui peut

Confidentialité

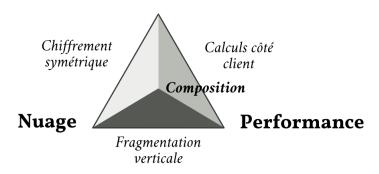


FIGURE 1 – Enjeux et techniques dans le cloud-computing

(image provenant de la thèse de Ronan Cherrueau)

améliorer les performances et permettre de tirer encore plus d'avantages du nuage.

La figure 1 résume comment ces trois techniques là interagissent avec les trois enjeux cités plus haut.

Dans sa thèse, Ronan Cherrueau montre qu'en composant ces différentes techniques on peut à la fois profiter des avantages du nuage, protéger la confidentialité et améliorer les performances d'un programme.

Pour pouvoir décrire comment s'effectue une telle composition, pour pouvoir vérifier la correction d'une telle composition et pour pouvoir raisonner dessus, Cherrueau a introduit un langage : C2QL.

2.2 Un langage pour décrire la composition : C2QL

Le langage C2QL donne une façon d'exprimer comment les techniques mentionnées cidessus se composent avec les fonctions classiques de l'algèbre relationnelle.

Les données manipulées sont donc représentées sous forme de tables, ou relations, c'est à dire un ensemble de lignes contenant des valeurs pour chacun(e) des différent(e)s attributs ou colonnes considéré(e)s.

Dans l'exemple de l'application stockant des rendez-vous, cette table contiendrait autant de lignes que de rendez-vous stockés dans l'application et (par exemple) trois colonnes ou attributs : nom de l'utilisateur ayant stocké le rendez-vous, date du rendez-vous et lieu du rendez-vous.

Les opérateurs empruntés à l'algèbre relationnelle présents dans ce langage sont

- La projection, notée π qui consiste à ne considérer que certains des attributs de la table.
- La sélection, notée σ qui consiste, pour une table donnée, à ne considérer que les lignes satisfiant un certain prédicat.
- La jonction naturelle, notée ⋈, qui combine les informations de deux tables en fonction des attributs qu'elles ont en commun.

— L'aggrégation et la réduction (notées respectivement group et fold) permettant, respectivement, de regrouper les lignes qui partagent les mêmes valeurs pour un certain ensemble d'attributs et d'effectuer des opérations sur les groupes de lignes ainsi obtenus.

Les opérateurs relatifs à la protection des données présents dans le langage sont

- La fragmentation verticale, notée frag, qui sépare une table en deux tables : la première contenant certains des attributs (colonnes) de la table d'origine, et la deuxième contenant le reste des attributs de la table d'origine.
- La défragmentation verticale, notée defrag, qui effectue l'opération inverse : à partir de deux tables n'ayant pas d'attributs en commun, reconstruit une seule table. Un attribut spécial, id, servant à identifier chaque ligne et créer lors de la fragmentation d'une table, rend la défragmentation possible.
- Le chiffrement, noté crypt qui chiffre un attribut dans une table.
- Le déchiffrement, noté decrypt qui déchiffre un attribut dans une table.

« Où est passée le calcul côté client? » est peut-être une question que vous vous posez peut-être en ce moment. La raison pour laquelle il n'y a aucun opérateur dans C2QL qui permette d'exprimer le rapatriement des données sur la machine de l'utilisateur est que cette gestion est faite implicitement, pour garantir que les requêtes C2QL protègent toujours la confidentialité, par conception. En effet, dans le langage C2QL on suppose que les données à protéger sont toujours soit un attribut, dans lequel cas il peut être protégé par chiffrement, soit une association de deux attributs, dans lequel cas elle peut être protégée par fragmentation. Par conséquent, dans une requête, le premier déchiffrement ou la première défragmentation (ou l'absence de chiffrement ou de fragmentation dans des cas où ils auraient été nécessaires) brisent la protection des données et requièrent que les données soient rapatriées chez l'utilisateur pour cette opération et toutes les opérations postérieures.

Ainsi, on peut, en regardant une expression C2QL et en connaissant les contraintes de confidentialité à respecter, déduire à quel moment les données doivent être rapatriées chez l'utilisateur; sans que cette opération aie besoin d'apparaître explicitement dans l'expression.

Rajouter un exemple serait probablement utile

TODO Expliquer les contraintes de sécurité avant l'aspect implicite du c.c.

Ce langage a été défini comme un Langage de Domaine Spécifique Embarqué dans Idris, qui est un langage à types dépendants. Le système de typage d'Idris et les types donnés aux opérateurs permettent de vérifier, lorsqu'on écrit une requête en C2QL que celle-ci aura un sens au moment de l'exécution.

2.3 Utiliser C2QL : commuter les opérateurs

Dans sa thèse, Cherrueau expose un méthode pour se servir facilement et efficacement de C2QL, en trois étapes.

Première étape : écrire la version locale de l'application. Le développeur peut commencer par écrire les requêtes de son application sans tenir compte de l'utilisation du cloud ni de la protection des données. C'est ce qu'on appelle ici une version « locale » du programme ; c'est

le programme tel qu'il pourrait s'exécuter dans la machine de l'utilisateur, localement. Cette version de l'application est beaucoup plus facile à écrire qu'une version utilisant efficacement les techniques de protection.

Deuxième étape : ajouter les techniques de protection nécessaires. Le langage C2QL suppose que les contraintes de confidentialité sont de deux types : soit il s'agit d'une association entre deux attributs que l'on veut protéger, dans lequel cas on choisit de mettre les deux attributs de l'association dans des fragments distincts grâce à la fragmentation verticale, soit il s'agit de la valeur d'un attribut que l'on veut protéger, dans lequel cas soit on chiffre cette valeur, soit on la décompose sur plusieurs fragments avec la technique exposée par Aggarwal (TODO rajouter la référence). Dans les deux cas on peut passer des versions locales des requêtes à des versions protégées en faisant apparaître à droite de la requête les fonctions de protections composées avec leurs fonctions réciproques.

Troisième étape : faire commuter les différentes fonctions. En faisant commuter les différents opérateurs qui apparaissent dans les requêtes, on peut aboutir à une requête optimisée, profitant des avantages du cloud et de bonnes performances. Pour cela, on a besoin de savoir à quelles conditions les différents opérateurs de C2QL peuvent commuter.

C'est sur cet ensemble de lois qui indiquent à quelle condition deux fonctions de C2QL peuvent commuter que je me suis intéressé pendant mon stage.

3 Contribution

Dans mon travail avec l'ensemble de lois de C2QL, j'ai commencé par poser une sémantique formelle pour les différentes fonctions du langage C2QL (section 3.1),puis j'ai complété cet ensemble de lois, en m'intéressé à toutes les combinaisons possibles de lois que l'on pouvait vouloir commuter dans le langage (section 3.2), ensuite j'ai démontré la correction sémantique de ces lois-là (section 3.3) et en ce moment je travaille sur l'automatisation de l'optimisation des requêtes, en réfléchissant à comment déterminer, parmi les commutations possibles, lesquelles il faut choisir pour optimiser une requête(section 3.4).

3.1 Établir des définitions

Les définitions des différents opérateurs présentes dans la thèse de 2016 sont données en français, ce qui facilite leur compréhension, mais rend impossible une preuve mathématique de la correction des lois les concernant.

J'ai donc commencé par poser des définitions formelles des fonctions du langage C2QL.

Pour former ces définitions, j'ai du faire des choix à plusieurs reprises. À chaque fois, j'ai utilisé les deux mêmes critères pour guider mes choix : d'une part il faut que les définitions que je choisi permettent au langage d'être le plus expressif possible, d'autre part il faut que les définitions choisies facilitent une démonstration rigoureuse de la correction des lois de commutation.

Des tuples ou des fonctions?

Dans son livre de 1982, Ullman parle de deux façons de définir les relations (tables) de l'algèbre relationnelle : soit comme un sous-ensemble d'un produit de domaines (donc comme un ensemble de tuples, où chaque tuple représente une ligne de la table et chaque coordonnée de chaque tuple correspond à la valeur pour un attribut donné), soit comme un ensemble de fonctions définies sur l'ensemble des attributs (chaque fonction correspond donc à une ligne de la table, et la valeur d'un attribut pour une ligne donnée est son image par cette fonction).

Ullman décide d'utiliser la première définition pour le reste de son livre, et ne détaille pas plus la deuxième définition. C'est pourtant cette deuxième définition que j'ai décidé de prendre dans ce cas ci, puisque que ce soit pour le chiffrement comme pour la fragmentation, on fait ici référence aux différentes colonnes d'une table par leur nom (toutes les colonnes ont un nom, et tous les noms des colonnes sont différents) et il est donc plus facile de raisonner, par exemple, sur la fragmentation, avec cette définition-là : en effet, la fragmentation se définit alors en thermes de simples restrictions sur certains ensembles.

Des sélections portant sur plus d'un attribut à la fois

- 3.2 Compléter les propriétés
- 3.3 Prouver les propriétés
- 3.4 Optimiser les requêtes
- 4 Travail futur
- 5 Conclusion