



**HAL**  
open science

## Opacité des artefacts d'un système workflow

Eric Badouel, Mohamadou Lamine Diouf

► **To cite this version:**

Eric Badouel, Mohamadou Lamine Diouf. Opacité des artefacts d'un système workflow. Revue Africaine de Recherche en Informatique et Mathématiques Appliquées, 2014, Volume 17 - 2014 - Special issue CARI'12, pp.117-196. 10.46298/arima.1973 . hal-00823005v2

**HAL Id: hal-00823005**

**<https://inria.hal.science/hal-00823005v2>**

Submitted on 8 Apr 2016

**HAL** is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.

## Opacité des artefacts d'un système workflow

Eric Badouel<sup>(a,c)</sup> et Mohamadou Lamine Diouf<sup>(b,c)</sup>

(a) Inria Rennes-Bretagne Atlantique, Campus Universitaire de Beaulieu,  
35042 Rennes Cedex, France

(b) Université Cheikh Anta Diop (UCAD), Dakar, Sénégal

(c) LIRIMA, Yaoundé.

{eric.badouel,mohamadou.diouf}@inria.fr

**RÉSUMÉ.** Une propriété d'un objet est dite *opaque* pour un observateur si celui-ci ne peut déduire que la propriété est satisfaite sur la base de l'observation qu'il a de cet objet. Supposons qu'un certain de nombre de propriétés (appelées *secrets*) soient attachées à chaque intervenant d'un système, nous dirons alors que le système lui-même est opaque si chaque secret d'un observateur lui est opaque: il ne peut percer aucun des secrets qui lui ont été attachés. L'opacité a été étudiée préalablement dans le contexte des systèmes à événements discrets où différents jeux d'hypothèses ont pu être identifiés pour lesquels on pouvait d'une part décider de l'opacité d'un système et d'autre part développer des techniques pour diagnostiquer et/ou forcer l'opacité. Ce papier constitue, à notre connaissance, la première contribution au problème de l'opacité des artefacts d'un système à flots de tâches (système workflow). Notre propos est par conséquent de formaliser ce problème en dégagant les hypothèses qui doivent être posées sur ces systèmes pour que l'opacité soit décidable. Nous indiquons quelques techniques pour assurer l'opacité d'un système.

**ABSTRACT.** A property (of an object) is opaque to an observer when he or she cannot deduce the property from its set of observations. If each observer is attached to a given set of properties (the so-called secrets), then the system is said to be opaque if each secret is opaque to the corresponding observer. Opacity has been studied in the context of discrete event dynamic systems where technique of control theory were designed to enforce opacity. To the best of our knowledge, this paper is the first attempt to formalize opacity of artifacts in data-centric workflow systems. We motivate this problem and give some assumptions that guarantee the decidability of opacity. Some techniques for enforcing opacity are indicated.

**MOTS-CLÉS :** Opacité, système à flots de tâches, artefact, documents structurés.

**KEYWORDS :** Opacity, data-centric workflow systems, artifact, structured documents.

## 1. Introduction

La propriété centrale qui nous intéresse dans ce travail est celle d'*opacité* [12, 5] dont nous donnons, dans le contexte de cet article, la présentation suivante. Soit un ensemble  $U$  —pour *univers*— contenant les objets qui nous intéressent. Dans le cas des systèmes à événements discrets il s'agit de l'ensemble des exécutions d'un système. Dans notre cas il s'agira de l'ensemble des dossiers manipulés par un système administratif. Ceux-ci sont représentés par les documents structurés —à la XML— conformes à un certain schéma (ou grammaire). Un observateur est représenté par sa fonction d'observation  $\phi : U \rightarrow O$  dans laquelle  $O$  est l'ensemble des observations possibles et  $\phi(x)$  est l'information que l'observateur a de l'objet  $x$ . Une propriété des objets peut être assimilée à l'ensemble des objets qui la vérifient, c'est-à-dire qu'il s'agit d'un ensemble  $P \subseteq U$  et on dit que " $x$  vérifie la propriété  $P$ " lorsque  $x \in P$  et que " $x$  ne vérifie pas la propriété  $P$ " si  $x \notin P$ . Une telle propriété  $P$  est *opaque* pour l'observateur si celui-ci ne peut pas déduire que l'objet qu'il observe satisfait la propriété sur la base de ce qu'il peut observer de celui-ci.

Formellement on peut exprimer que  $P$  est alors opaque vis-à-vis de  $\phi$  lorsque

$$\forall x \in P \exists y \notin P \phi(x) = \phi(y) \quad (1)$$

c'est-à-dire qu'il n'est pas possible de déduire qu'un objet vérifie la propriété ( $x \in P$ ) sur la base de son observation  $\phi(x)$  puisqu'il est en effet possible de trouver un autre élément  $y$  donnant lieu à la même observation ( $\phi(x) = \phi(y)$ ) et qui néanmoins ne vérifie pas cette propriété ( $y \notin P$ ). L'équation (1) peut se reformuler de façon ensembliste : la propriété  $P$  est opaque vis-à-vis de  $\phi$  si, et seulement si,

$$\phi(P) \subseteq \phi(U \setminus P) \quad (2)$$

L'hypothèse sous-jacente à la notion d'opacité est que l'observateur a une parfaite connaissance du système : il connaît l'ensemble  $U$  —dans notre cas la grammaire décrivant les documents—, il sait par ailleurs comment sa fonction d'observation  $\phi$  est construite et il est en mesure de déterminer l'ensemble  $\phi^{-1}(o)$  des objets qui donnent lieu à une observation donnée  $o \in O$ . De cette façon, la connaissance qu'il a d'un objet  $x \in U$  est donnée par l'ensemble  $\phi^{-1}(\phi(x))$  des objets qui donnent lieu à la même observation que celui-ci. Ainsi il sera capable de déterminer qu'un objet  $x$  vérifie la propriété  $P$  lorsque

$$\phi^{-1}(\phi(x)) \subseteq P \quad (3)$$

De fait les formules —équivalentes— (1) et (2) peuvent se reformuler en disant que la condition (3) ne doit être vérifiée par aucun élément  $x \in U$ . Lorsque la relation (3) est vérifiée pour un certain objet  $x$  on dit que cet objet donne lieu à une *fuite d'information*, ce qui constitue une faille de sécurité dans le cas où la propriété qui a été déduite est une information sensible qu'on souhaitait garder secrète. L'élément  $x$  vérifiant la relation (3) est appelé un *témoin* de la fuite de l'information  $P$ . Il faut faire attention au fait que la non opacité d'une propriété  $P$  ne signifie pas qu'un observateur sera toujours dans la possibilité de déterminer si la propriété  $P$  est vraie ou non pour l'objet qu'il observe mais seulement qu'il existe des cas —donnés par les témoins— pour lesquels il est capable de déduire que la propriété est satisfaite. On doit noter par ailleurs que l'opacité d'une propriété  $P$  n'est pas équivalente à l'opacité de sa négation : il peut exister des

témoins pour  $P$  sans qu'il existe des témoins pour  $\neg P$ . Cette asymétrie de l'opacité est particulièrement visible dans la formulation (2) de cette propriété.

Supposons qu'à chaque intervenant du système (identifié à sa fonction d'observation  $\phi : U \rightarrow O$ ) soit attaché un certain nombre de propriétés  $S_1, \dots, S_n \subseteq U$ . Il s'agit de propriétés du système, appelées *secrets*, qui ne doivent pas être dévoilées à cet utilisateur. Typiquement certaines informations confidentielles doivent être accessibles à certains utilisateurs autorisés et être tenues secrètes pour les autres. Par exemple les coordonnées bancaires d'un client doivent être connues du service de comptabilité mais pas des autres services du système. Nous exigeons alors qu'aucun secret ne puisse être dévoilé : le système est *opaque* si pour chaque observateur  $\phi : U \rightarrow O$  et pour chacun des secrets  $S \subseteq U$  attachés à cet observateur on a  $\phi(S) \subseteq \phi(U \setminus S)$ , c'est-à-dire que le secret  $S$  est opaque vis-à-vis de  $\phi$ .

Le problème de l'opacité a été jusqu'à présent étudié dans le contexte des systèmes à événements discrets et a servi à modéliser des problèmes de sécurité et/ou de confidentialité pour des systèmes informatiques ou des protocoles, voir e.g. [5, 13]. Nous revenons sur ces travaux dans la conclusion lorsqu'il s'agira d'envisager des techniques pour assurer l'opacité d'un système en s'inspirant des techniques de monitoring [9] ou de contrôle [1] qui ont été développées dans ce contexte d'origine.

Notre travail est la première tentative, à notre connaissance, d'appliquer ce concept d'opacité à des systèmes à flots de tâches centrés sur les données (*data-centric workflow systems* [7]). Il s'agit de systèmes distribués asynchrones dans lesquels la coordination des activités s'effectue par la transmission de documents structurés (à la XML) combinant structure logique, et données. Les éléments de l'univers de référence  $U$  sont donc des documents structurés, appelés *artefacts*, qui servent de support à la réalisation d'une tâche : il s'agira par exemple du dossier médical d'un patient, d'un dossier administratif, d'un dossier pour le suivi d'une commande ... Au cours de ce traitement nous avons besoin d'invoquer un certain nombre de services auxquels nous serons ainsi amené à transmettre des informations extraites de ce dossier. Il y a plusieurs raisons qui conduisent à ne transmettre qu'une partie des informations contenues dans l'artefact. D'une part nous voulons éviter de surcharger un service par des informations qui ne sont pas utiles à la réalisation de la tâche qui lui incombe. D'autre part pour des raisons de confidentialité, et c'est ce point qui nous intéresse ici, certaines informations sensibles ne doivent pas être connues de tous. Enfin les services invoqués ont généralement été définis antérieurement et de façon par conséquent indépendante du workflow qui les invoque ; celui-ci doit donc extraire de l'artefact un document conforme à ce que le service est préparé à recevoir.

La discussion précédente va nous guider dans le choix des structures grammaticales (statuant de la conformité d'un document) ainsi que sur les fonctions d'observations  $\phi : U \rightarrow O$  utilisées pour l'extraction des informations à transmettre à un service donné. Un autre critère important sur le choix de ces deux paramètres est de pouvoir aboutir à des algorithmes permettant de vérifier de manière effective si un système donné est opaque ou non. Un de nos principaux objectifs est effectivement d'aboutir à une formalisation des notions de structure grammaticale des documents, de fonctions d'observation et de secrets qui garantisse la décision de l'opacité. Concrètement nous souhaitons obtenir un algorithme nous indiquant si un système est opaque ou non et qui dans la négative nous construit des témoins pour chacun des secrets qui peuvent être dévoilés.

Dans la conclusion nous donnons quelques éléments sur la façon dont on peut modifier le système pour assurer son opacité. La première approche, qui s'inspire des techniques

de contrôle des systèmes à événements discrets, permet de construire une nouvelle grammaire qui restreint les documents de la grammaire d'origine à ceux pour lesquels aucun secret n'est dévoilé. Cette technique peut avoir l'inconvénient d'écarter des documents que le concepteur du système aurait par ailleurs toutes les raisons de vouloir considérer comme parfaitement licites. Une méthode moins intrusive, qui s'inspire de techniques d'extensions conservatives de schémas XML, consiste au contraire à étendre la grammaire.

---

## 2. Conformité d'un artefact

Les arbres construits sur un alphabet  $\Omega$  —c'est-à-dire un ensemble fini— sont définis inductivement comme suit

$$t := \omega(t_1, \dots, t_n) \quad \omega \in \Omega$$

$\omega$  est l'étiquette attachée à la racine de l'arbre et  $t_1, \dots, t_n$  est la liste des sous arbres se trouvant sous le noeud racine. Cette liste peut-être vide ( $n = 0$ ) auquel cas l'arbre  $t = \omega()$  est réduit à son étiquette —on pourra écrire de façon abrégée  $t = \omega$  dans ce cas—. Une forêt est une liste d'arbres.

Un artefact est un ensemble de documents constituant le dossier pour un client d'un certain service. Ces documents ont une structure d'arbre. Les informations attachées aux noeuds d'un arbre sont données par des paires attributs/valeurs qui permettent de collecter les informations pertinentes utiles pour le traitement du dossier. Néanmoins comme on aura besoin de travailler avec un ensemble fini d'étiquettes on suppose donnée une abstraction finie de ces informations sous la forme d'un alphabet  $\Omega$  utilisé pour étiqueter les différents noeuds des documents. La structure arborescente reflète l'organisation logique du document :  $d = \omega(d_1, \dots, d_n)$  signifie que  $d$  est un document portant une information  $\omega$  et dont  $d_1, \dots, d_n$  sont les constituants (les sous dossiers). Comme une forêt  $f = d_1 \dots d_n$  peut-être encodée par l'expression  $\#(d_1, \dots, d_n)$  et quitte à rajouter le symbole  $\#$  à l'ensemble  $\Omega$ , nous pourrions supposer qu'un artefact est également donné par un arbre.

Dans cette section nous nous intéressons à caractériser la conformité d'un arbre selon une grammaire (d'arbres). Cela permet de fixer la structure logique du document, c'est-à-dire la façon dont ses différents constituants sont organisés. La notion de grammaire considérée ici correspond aux automates d'arbres (et de forêts) introduits par Murata et al [4]. Il s'agit d'une définition assez générale dans laquelle l'arité d'un noeud n'est pas fixe —comme cela est le cas pour une grammaire ordinaire, c'est-à-dire pour une signature multi sortes—. Ceci est utile lorsque parmi des constituants d'un sous dossier on peut trouver une liste d'éléments d'une certaine catégorie (par exemple la liste des ouvrages d'une bibliothèque) et dont le nombre n'est par conséquent pas déterminé, ou lorsque certains éléments sont optionnels. Le plus simple est alors d'utiliser une expression régulière pour décrire ce qu'on est censé y trouver par exemple la règle

$$A \rightarrow \omega_1 \langle B^*(A + \varepsilon) \rangle + \omega_2 \langle \varepsilon \rangle$$

s'interprète de la manière suivante.  $A$  et  $B$  sont des sortes (aussi appelées catégories syntaxiques). L'ensemble des sortes est fini et sert à cataloguer les différents constituants d'un dossier en un nombre (donc fini) de catégories, par exemple ceci est un RIB, cet

autre élément est une adresse ou un dossier de réservation d'hôtel. Comme dit plus haut les étiquettes  $\omega \in \Omega$  sont des informations extraites des données se trouvant en racine du dossier (afin de travailler ici aussi avec un ensemble fini de possibilités). La règle précédente peut se lire alors comme suit. Un arbre est bien formé et est de la sorte  $A$  si l'information extraite à sa racine est  $\omega_1$  et il possède comme fils une suite d'éléments de sorte  $B$  —dont le nombre  $n$  est pas fixé— et possiblement un élément de sorte  $A$ , ou bien l'information extraite à sa racine est  $\omega_2$  et dans ce cas il ne doit avoir aucun successeur.

Avant de définir formellement les grammaires nous introduisons la notion plus générale d'algèbre qui est mathématiquement plus simple et donc plus adaptée pour énoncer et prouver des résultats techniques. Une algèbre est un mécanisme permettant d'extraire (ou de synthétiser) des informations à partir du document par induction sur sa structure. Exécuter une requête sur le document pourra ainsi se faire en évaluant le document selon une algèbre. Si l'information ainsi extraite est une valeur de vérité (la valeur retournée par l'évaluation est booléenne) alors l'algèbre peut être vue comme une propriété des documents. En particulier un secret, vue comme une propriété des artefacts, sera identifié à une algèbre particulière.

**Définition 2.1** (Algèbre). *Une  $\Omega$ -algèbre de support  $D$  est une application  $\mathcal{A} : \Omega \times D^* \rightarrow D$ . On note  $\omega^A : D^* \rightarrow D$ , appelée **fonction d'interprétation** de  $\omega \in \Omega$ , l'application définie par :  $\omega^A(x_1, \dots, x_n) = \mathcal{A}(\omega, x_1, \dots, x_n)$ . L'ensemble  $T(\Omega)$  des arbres étiquetés sur  $\Omega$  est le plus petit ensemble tel que*

$$(t_1, \dots, t_n \in T(\Omega) \wedge \omega \in \Omega) \implies \omega(t_1, \dots, t_n) \in T(\Omega)$$

*La valeur  $t^A$  d'un arbre  $t \in T(\Omega)$  selon l'algèbre  $\mathcal{A}$  est définie inductivement par :*

$$t = \omega(t_1, \dots, t_n) \implies t^A = \omega^A(t_1^A, \dots, t_n^A)$$

*Si  $\mathcal{A}$  est une algèbre et  $x \in D$  un élément de son support on pose*

$$L(\mathcal{A}, x) = \{t \in T(\Omega) \mid t^A = x\}$$

NOTATIONS. — Notons  $\mathcal{A}_{\omega, x} = (\omega^A)^{-1} = \{u \in D^* \mid \mathcal{A}(\omega, u) = x\}$  pour  $\omega \in \Omega$  et  $x \in D$ . Pour un symbole  $\omega$  donné ces ensembles sont donc disjoints :  $x \neq y \implies \mathcal{A}_{\omega, x} \cap \mathcal{A}_{\omega, y} = \emptyset$ . On note  $\mathcal{A}_\omega = \bigcup_{x \in D} \mathcal{A}_{\omega, x}$  leur union. De façon similaire nous associons à  $x \in D$  l'ensemble  $\mathcal{A}_x = \bigcup_{\omega \in \Omega} \mathcal{A}_{\omega, x}$ .

Une algèbre est non déterministe si le résultat de l'évaluation peut produire différentes valeurs. C'est-à-dire qu'il s'agit d'une fonction  $\mathcal{A} : \Omega \times D^* \rightarrow \wp(D)$ . On peut alors de façon classique se ramener au cas déterministe en passant aux parties c'est-à-dire en remplaçant le support  $D$  de l'algèbre par l'ensemble  $\wp(D)$  des parties de  $D$ .

REMARQUE. — **Algèbre non-déterministe** — Une application  $\mathcal{A} : \Omega \times D^* \rightarrow \wp(D)$ , appelée algèbre non-déterministe, est assimilée à l'algèbre  $\mathcal{A}^d : \Omega \times \wp(D)^* \rightarrow \wp(D)$  définie sur  $\wp(D)$  par

$$\mathcal{A}^d(\omega, X_1 \cdots X_n) = \bigcup \{\mathcal{A}(\omega, x_1 \cdots x_n) \mid \forall i \in \{1, \dots, n\} \ x_i \in X_i\}$$

Le langage d'une algèbre non déterministe est ainsi donné par  $L(\mathcal{A}, x) = L(\mathcal{A}^d, x)$ , c'est-à-dire  $L(\mathcal{A}, x) = \{t \in T(\Omega) \mid x \in t^A\}$ .

La première propriété d'un artefact qu'il convient de vérifier est sa conformité à un certain schéma donné par une grammaire. Comme mentionné plus haut la notion de grammaire utilisée ici s'inspire des automates d'arbres de [4].

**Définition 2.2** (Grammaire). *Une grammaire  $G = (\Omega, \Xi, \mathcal{L})$  sur  $\Omega$  est la donnée d'un ensemble fini  $\Xi$  de sortes et d'un langage régulier  $\mathcal{L}_{\omega, s} \subseteq \Xi^*$  associé à chaque  $\omega \in \Omega$  et  $s \in \Xi$ . On note  $G \vdash t :: s$  pour signifier que  $t$  est un arbre conforme à la grammaire  $G$  et est de sorte  $s$ . On note  $L(G, s) = \{t \mid G \vdash t :: s\}$  cet ensemble. Ces ensembles sont définis inductivement comme suit :*

$$(\forall i \in \{1, \dots, n\} \ G \vdash t_i :: s_i \ \wedge \ s_1 \cdots s_n \in \mathcal{L}_{\omega, s}) \Rightarrow G \vdash \omega(t_1, \dots, t_n) :: s$$

Le langage régulier  $\mathcal{L}_{\omega, s} \subseteq \Xi^*$  décrit les sortes possibles des constituants (noeuds fils) d'un noeud de sorte  $s$ . Par exemple  $\mathcal{L}_{\omega, A} = B^* + C$  exprime le fait qu'un noeud  $\omega$  sera de sorte  $A$  si ses successeurs immédiats sont soit une liste (éventuellement vide) de noeuds de sorte  $B$  soit un unique noeud de sorte  $C$ . Remarquons que ces noeuds ont une arité variable : l'étiquette d'un noeud ne caractérise pas le nombre de ses successeurs.

NOTATIONS. — **Grammaire sous forme de productions** — Une grammaire sera présentée syntaxiquement sous la forme d'un ensemble de **productions**  $A \rightarrow \omega \langle E \rangle$  où  $E$  est une expression régulière telle que  $L(E) = \mathcal{L}_{\omega, A}$  lorsque cet ensemble est non vide.

La vérification de la conformité d'un arbre peut-être implémentée par un automate à pile. Cet automate à pile sera déterministe si  $s \neq s' \Rightarrow \mathcal{L}_{\omega, s} \cap \mathcal{L}_{\omega, s'} = \emptyset$ , c'est-à-dire que la sorte d'un noeud est caractérisée par son étiquette et la sorte de ses successeurs. Lorsque cette propriété est vérifiée la grammaire est dite *déterministe*. Comme dans le cas des algèbres, une grammaire non-déterministe peut être transformée en une grammaire déterministe équivalente (c'est-à-dire ayant le même langage).

**Définition 2.3** (Déterminisation d'une grammaire). *La déterminisation d'une grammaire  $G = (\Omega, \Xi, \mathcal{L})$  est la grammaire déterministe  $G^d = (\Omega, \wp(\Xi) \setminus \{\emptyset\}, \mathcal{L}^d)$  pour laquelle  $\mathcal{L}_{\omega, X}^d = \bigcup_{s \in X} \mathcal{L}_{\omega, s}$ .*

La remarque suivante montre que les grammaires ne sont rien d'autres que les algèbres régulières, c'est-à-dire les algèbres dont les langages  $\mathcal{A}_{\omega, x}$  sont réguliers. Par défaut néanmoins une algèbre est déterministe alors qu'une grammaire est par défaut non-déterministe.

REMARQUE. — **Grammaires vs algèbres** – Une grammaire  $G$  peut-être assimilée à l'algèbre non-déterministe  $\mathcal{A}$  de domaine  $\Xi$  telle que  $\mathcal{A}(\omega, u) = \{s \in \Xi \mid u \in \mathcal{L}_{\omega, s}\}$  pour  $\omega \in \Omega$ ,  $s \in \Xi$  et  $u \in \Xi^*$ , et donc  $\mathcal{A}_{\omega, s} = \mathcal{L}_{\omega, s}$ . Une grammaire déterministe  $G$  peut être assimilée à l'algèbre de domaine  $D = \Xi \cup \{\top\}$  où  $\top$  est un symbole supplémentaire ( $\top \notin \Xi$ ) et

$$\omega^G(s_1, \dots, s_n) = \begin{cases} s & \text{si } s_1 \cdots s_n \in \mathcal{L}_{\omega, s} \\ \top & \text{sinon} \end{cases}$$

Les opérations respectives de déterminisation des grammaires et des algèbres se correspondent, c'est-à-dire que  $G^d$  décrite dans la définition 2.3 est la déterminisation de  $G$  vue comme algèbre non-déterministe : l'élément additionnel  $\top$  correspond à la partie vide qu'on avait pris soin de ne pas introduire dans la définition de  $G^d$ . La correspondance

entre une grammaire déterministe et l'algèbre qui lui est associée est donnée par l'identité  $L(G, s) = \{t \in T(\Omega) \mid t^G = s\}$  qui découle immédiatement des définitions 2.1 et 2.2. Un arbre est ainsi conforme à la grammaire (ou bien formé) si et seulement si  $t^G \neq \top$ . Remarquons que comme  $\mathcal{L}_{\omega, s} \subseteq \Xi^*$  et  $\top \notin \Xi$  il vient que  $t = \omega(t_1, \dots, t_n)$  est non conforme (i.e.,  $\omega^G(t_1^G, \dots, t_n^G) = \top$ ) dès qu'un de ses sous-arbres est non conforme.

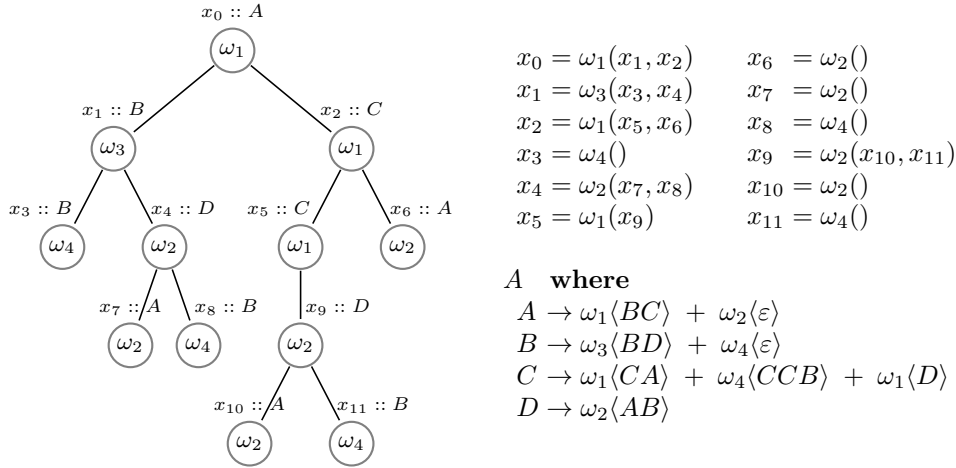
NOTE. — **Langage d'une grammaire** – Dans la suite nous ferons l'hypothèse suivante sur les grammaires. L'ensemble des sortes contient un symbole spécifique  $\mathbf{ax} \in \Xi$ , appelé **axiome** de la grammaire, pour lequel il est associé une unique production  $\mathbf{ax} \rightarrow \#(E_0)$  telle que le symbole  $\# \in \Omega$  (appelé racine) n'apparaît dans aucune des autres productions de la grammaire.  $t^G = s$  pour  $t \in T(\Omega \setminus \{\#\})$  et  $s \in \Xi \setminus \{\top\}$  signifie que  $t$  est bien formé de sorte  $s$ , et on pose  $L(G) = \{t_1 \dots t_n \in T(\Omega \setminus \{\#\})^* \mid t_1^G, \dots, t_n^G \in L(E_0)\}$  le langage de la grammaire  $G$ , i.e.,  $t_1 \dots t_n \in L(G)$  si, et seulement si  $(\#(t_1, \dots, t_n))^G = \mathbf{ax}$ .

En utilisant les résultats classiques sur les langages reconnaissable d'arbres [4, 8] on peut montrer que les langages des grammaires sont clos de manière effective par les opérations booléennes et qu'on peut décider de la vacuité et donc de l'inclusion de ces langages, ainsi la relation  $G \leq G' \iff L(G) \subseteq L(G')$  est décidable.

Dans la remarque suivante —un peu longue mais pouvant être omise par le lecteur non intéressé— nous faisons une digression qui peut permettre au lecteur de mieux comprendre le type de propriétés d'arbres (en particulier les secrets) qu'on souhaite pouvoir exprimer en utilisant des algèbres. Cet aspect ne sera pas développé dans la suite de cet article.

REMARQUE. — Dans ce travail préliminaire nous ne cherchons pas à identifier une syntaxe pour exprimer les secrets sur les documents. De façon simplifiée nous avons dit précédemment qu'on pouvait assimiler un secret (en tant que propriété des documents) à une algèbre dont le résultat retourné en racine de l'arbre est booléen (vrai ou faux). Cela n'exclut pas que le support de l'algèbre contienne d'autres valeurs de nature bien plus complexe afin de représenter des valeurs intermédiaires associées à des sous-parties du document. Dans la mesure où les informations attachées à un noeud sont des paires attributs/valeurs un formalisme particulièrement adapté pour exprimer ces calculs est celui des grammaires attribuées. Il est bien connu que l'évaluation des attributs pour une grammaire attribuée peut s'exprimer par une algèbre, néanmoins d'ordre supérieur, dont les fonctions d'interprétation permettent en chaque noeud de l'arbre de calculer les attributs synthétisés en fonction des attributs hérités. Il est montré dans [3] que les règles sémantiques d'une grammaire attribuée peuvent se décrire par une algèbre *du premier ordre* sur une signature étendue dans laquelle aux liens d'un père vers ses fils on rajoute un lien supplémentaire de ce noeud vers son propre père afin d'accéder à son contexte. Pour des raisons d'effectivité on doit se restreindre ici à des algèbres dont le support est fini et donc considérer des abstractions des informations attachées aux noeuds et des règles sémantiques qui permettent de les calculer. Ce sera plus facile à obtenir si nous partons d'une algèbre du premier ordre. D'où l'intérêt ici de pouvoir manipuler des liens symboliques permettant d'exprimer des propriétés des documents (en particulier des secrets) qui ne seront pas nécessairement définis par induction sur leur structure physique mais pourront aussi dépendre d'informations contextuelles.





**Figure 1.** Un arbre vu comme un système d'équations dont la résolution dans une algèbre (par exemple une grammaire déterministe) donne la valeur de l'arbre. Le préambule "L where ..." à la liste des productions de la grammaire représente la production  $ax \rightarrow \sharp\langle L \rangle$  associée à son axiome (qui n'est, ainsi que le symbole en racine  $\sharp$ , pas représenté explicitement).

Pour ce faire, et comme illustré à la figure 1, nous pouvons décrire un arbre par un système d'équations  $E$  de la forme  $x = \omega(x_1, \dots, x_n)$  dans laquelle la variable  $x$  est dite *définie* et les variables  $x_i$  ( $1 \leq i \leq n$ ) sont dites *utilisées*. Les variables correspondent bijectivement aux noeuds de l'arbre. Chacune d'elles doit être définie par exactement une équation et ne doit être utilisée qu'au plus une fois. Enfin la relation de dépendance, qui est donnée par  $x \ll y$  ssi " $y$  est utilisée dans l'équation qui définit  $x$ ", doit être acyclique. L'évaluation de l'arbre dans une algèbre revient à associer une valeur  $x^A$  pour chaque variable de sorte que  $x = \omega(x_1, \dots, x_n) \implies x^A = \omega^A(x_1^A, \dots, x_n^A)$ . Autrement dit le vecteur  $\vec{x}^A$  est solution du système  $E_A$ , dit *interprétation de  $E$  dans  $A$* , formé des équations  $x = \omega^A(x_1, \dots, x_n)$ , c.-à-d.  $\vec{x}^A = E_A(\vec{x}^A)$ . Les composantes de ce vecteur  $\vec{x}^A : X \rightarrow D$  peuvent être calculées suivant un tri topologique du graphe de dépendance (on évalue un noeud après évaluation des noeuds dont il dépend). L'avantage de cette présentation équationnelle de l'évaluation d'un document (par opposition à une présentation plus opérationnelle à l'aide de machines à piles) est de pouvoir s'appliquer sans modification au cas où on souhaite pouvoir faire du partage de documents : l'hypothèse qu'une variable n'est jamais utilisée de façon multiple sert à caractériser les arbres mais n'a aucune utilité dans la résolution des systèmes d'équations correspondants.

On peut également lever l'hypothèse que le graphe de dépendance est acyclique en résolvant le système d'équation  $\vec{x}^A = E_A(\vec{x}^A)$  de manière itérative. Il faut pour cela supposer que le domaine de l'algèbre est un treillis complet (dont le plus petit élément  $\perp$ , la valeur indéfinie, est la valeur partagée par tout objet et le plus grand élément  $\top$ , la valeur contradictoire, est une valeur qui ne peut-être prise par aucun objet) et que les fonctions d'interprétations sont continues (c'est-à-dire commutent aux bornes supérieures). Dans ce cas le plus petit point fixe de l'équation  $\vec{x} = E_A(\vec{x})$ , noté  $E_A^\dagger$  est la borne supérieure de la suite obtenue en itérant le constructeur  $E_A$  à partir du vecteur totalement indéfini  $\vec{\perp}$  (dont toutes les composantes sont égales à  $\perp$ ), c'est-à-dire que  $E_A^\dagger \vee E_A^n(\vec{\perp})$ . Cer-

taines hypothèses sur le support de l'algèbre (comme le fait qu'il soit de hauteur bornée) garantissent que la suite croissante  $\langle E_{\mathcal{A}}^n(\bar{\perp}) \rangle_{n \in \mathbb{N}}$  est stationnaire, c'est-à-dire se stabilise toujours à partir d'un rang  $m$  pour lequel  $E_{\mathcal{A}}^{\dagger} = E_{\mathcal{A}}^m(\bar{\perp})$ .

Un cycle dans la relation de dépendance semble dire qu'un document est une composante de lui-même, ce qui ne fait pas de sens. Néanmoins dans certains cas la relation de dépendance ne correspondra pas uniquement à la relation de contenance (la structure physique sous-jacente du document) mais aussi à des liens symboliques qui peuvent être utilisés pour rechercher des informations. C'est-à-dire qu'on trouvera un chemin dans le graphe de dépendance de  $x$  vers  $y$  si le calcul d'un certain attribut du noeud  $x$  peut dépendre de la valeur de certains attributs du noeud  $y$ . Même si on ne doit pas trouver de cycle dans les dépendances entre les attributs on pourra en trouver entre les noeuds qui portent ces attributs dès que la valeur d'un attribut peut dépendre d'informations provenant du contexte de ce noeud (et non seulement de ses sous-arbres).

**REMARQUE. — Grammaires vs signatures multi-sortes** – Les grammaires sont une généralisation des signatures multisortes qui peuvent-être assimilées aux grammaires dans lesquelles pour chaque symbole  $\omega \in \Omega$  les langages  $\mathcal{L}_{\omega, s}$  sont tous vides sauf un qui est réduit à un seul mot. C'est-à-dire que chaque symbole a une arité et une sorte qui sont fixées : on note  $\omega :: s_1 \dots s_n \rightarrow s$  lorsque  $\mathcal{L}_{\omega, s} = \{s_1 \dots s_n\}$  (et donc  $s \neq s' \implies \mathcal{L}_{\omega, s'} = \emptyset$ ). Inversement, à toute grammaire déterministe on peut associer une signature multisorte dont les opérateurs sont les paires  $\langle \omega, u \rangle$  pour  $u = s_1 \dots s_n \in \Xi^*$  telles que  $G(\omega, u) = s \neq \top$  il s'agit alors d'un opérateur d'arité  $u$  et de sorte  $s$ , c'est-à-dire  $\langle \omega, u \rangle :: s_1 \dots s_n \rightarrow s$ . Rappelons qu'une algèbre  $\mathcal{A}$  pour une signature multi-sortes est donnée par un ensemble non vide  $\mathcal{A}_s$  associé à chaque sorte  $s \in \Xi$  et une fonction d'interprétation  $\omega^{\mathcal{A}} : \mathcal{A}_{s_1} \dots \mathcal{A}_{s_n} \rightarrow \mathcal{A}_s$  associée à chaque opérateur  $\omega :: s_1 \dots s_n \rightarrow s$ . Les algèbres pour la signature multisorte associée à une grammaire déterministe  $G$  sont les  $\Omega$ -algèbres  $\mathcal{A} : \Omega \times D^* \rightarrow D$  pour lesquelles il existe une application surjective  $\varphi : D \rightarrow \Xi$  telle que  $\mathcal{L}_{\omega, \varphi(x)} = \{\varphi(x_1) \dots \varphi(x_n) \mid \mathcal{A}(\omega, x_1 \dots x_n) = x\}$ . Autrement dit  $G \circ \langle id_{\Omega}, \varphi^* \rangle = \varphi \circ \mathcal{A}$ . Nous dirons par analogie que ces algèbres sont les **algèbres de la grammaire**.

Un artefact (liste de documents) est conforme s'il appartient au langage de la grammaire et en particulier chacun des documents qui le constitue est bien formé. Une fois que la conformité de l'artefact est acquise on peut en extraire de l'information en l'évaluant selon des algèbres de la grammaire au sens de la remarque précédente.

---

### 3. Abstraction d'un artefact

Dans la section précédente nous nous sommes intéressés à la structure logique d'un artefact (donnée par une grammaire) et à son évaluation selon une algèbre. Nous allons maintenant nous intéresser aux données qu'il contient. Les informations attachées à un noeud correspondent à son étiquette  $\omega \in \Omega$ . Nous pouvons voir un artefact comme des boîtes imbriquées avec pour chacune d'entre elles les informations qui lui sont attachées. Une abstraction (ou observation) consiste à supprimer toutes les boîtes associées à des sortes "non visibles" en supprimant par la même occasion les données qui leur correspondent. L'effet de cette transformation sur la structure du document peut se décrire comme suit.

**Définition 3.1.** La fonction de projection  $p_{\Xi'} : T(\Omega) \rightarrow T(\Omega)^*$  associée à un sous-ensemble  $\Xi' \subseteq \Xi$  (les sortes visibles) est donnée pour  $t = \omega(t_1, \dots, t_n)$  par :

$$p_{\Xi'}(t) = \begin{cases} \omega(p_{\Xi'}(t_1) \dots p_{\Xi'}(t_n)) & \text{si } t^G \in \Xi' \\ p_{\Xi'}(t_1) \dots p_{\Xi'}(t_n) & \text{si } t^G \notin \Xi' \end{cases}$$

Les sortes visibles contiennent toujours le symbole en racine :  $\# \in \Xi'$ .

Ainsi si la sorte d'un arbre est visible sa projection est un arbre sinon il s'agit d'une suite d'arbres —i.e., une forêt—.

**Définition 3.2.** Si  $\pi = p_{\Xi'}$  est la projection associée à  $\Xi' \subseteq \Xi$ , on pose  $\mathcal{A} = G/\pi$  l'algèbre non-déterministe  $\mathcal{A} : \Omega \times \Xi'^* \rightarrow \wp(\Xi')$  définie par

$$\mathcal{A}(\omega, v) = \{s \in \Xi' \mid \exists u \in \mathcal{L}_{\omega, s} \quad u \vdash^* v\}$$

où  $\vdash^*$  est la fermeture réflexive et transitive de la relation

$$\vdash = \{(u \cdot s \cdot u', u \cdot v \cdot u' \mid s \in \Xi \setminus \Xi' \text{ et } v \in \mathcal{L}_s)\}$$

Ainsi  $\mathcal{A}^d(\omega, X_1 \dots X_n) = \{s \in \Xi' \mid \exists u \in \mathcal{L}_{\omega, s} \exists x_i \in X_i \quad u \vdash^* x_1 \dots x_n\}$ .

REMARQUE. — **Dérivation maximale** – On note  $u \vdash_{\max}^* v$  lorsque  $u \vdash^* v$  et cette dérivation est maximale, c'est-à-dire qu'il n'existe pas de  $w$  dans lequel  $v$  se dérive. Si les langages  $\mathcal{L}_s = \cup_{\omega} \mathcal{L}_{\omega, s}$  sont tous non vides —hypothèse que nous ferons par la suite— alors

$$u \vdash_{\max}^* v \iff u \vdash^* v \text{ et } v \in (\Xi')^*$$

Le théorème suivant est le principal résultat de cet article. Il montre la correction de la grammaire quotient en ce sens qu'elle reconnaît exactement les projections des arbres conformes à la grammaire de départ.

**Théorème 3.3.**  $L(G/\pi) = \pi(L(G))$ .

Afin d'établir ce théorème nous introduisons sur les arbres une relation, indexée par les contextes, qui correspond à la relation de dérivation  $\vdash$  sur les sortes.

**Définition 3.4** (Contextes). Les opérateurs dérivés, ou **contextes**,  $C \in \Omega^{\textcircled{n}}$  sont les arbres construits sur la signature  $\Omega \cup \{[\ ]\}$  obtenue à partir de  $\Omega$  en ajoutant un symbole spécial  $[\ ]$  figurant la position d'un argument du contexte. Les contextes sont ainsi définis par la syntaxe BNF suivante :

$$C ::= [\ ] \mid \omega(C_1, \dots, C_n)$$

L'**arité** d'un contexte, correspondant à son nombre d'arguments, est défini inductivement par  $\text{ar}([\ ]) = 1$  et  $\text{ar}(\omega(C_1, \dots, C_n)) = \text{ar}(C_1) + \dots + \text{ar}(C_n)$  (avec par convention  $\text{ar}(\omega()) = 0$  correspondant au cas où  $n = 0$ ). On note  $\Omega_n^{\textcircled{n}}$  l'ensemble des contextes d'arité  $n$  et  $\omega_n \in \Omega_n^{\textcircled{n}}$  le contexte d'arité  $n$  associé à l'opérateur  $\omega$ , il s'agit du contexte  $\omega_n = \omega([\ ], \dots, [\ ])$  contenant exactement  $n$  arguments. Un contexte  $C \in \Omega_n^{\textcircled{n}}$  d'arité  $n$  s'interprète selon une algèbre  $\mathcal{A} : \Omega \times D^* \rightarrow D$  comme l'opérateur  $C^{\mathcal{A}} : D^n \rightarrow D$  donné par :

- (1)  $[\ ]^A(x) = x$
- (2)  $(\omega(C_1, \dots, C_n))^A(v) = \omega^A(C_1^A(v_1), \dots, C_n^A(v_n))$   
où  $|v_i| = \mathbf{ar}(C_i)$  et  $v_1 \cdots v_n = v$ .

En particulier on assimilera un contexte avec son interprétation dans l'algèbre libre  $F : \Omega \times T(\Omega)^* \rightarrow T(\Omega)$  (pour laquelle  $\omega^F(t_1, \dots, t_n) = \omega(t_1, \dots, t_n)$ ), c'est-à-dire que l'expression  $C(t_1, \dots, t_n)$  sera définie si, et seulement si,  $\omega$  est d'arité  $n$ , et dans ce cas elle correspond à l'arbre obtenu en substituant les arbres  $t_1, \dots, t_n$  aux arguments du contexte pris dans leur ordre d'apparition. L'écriture  $t = C(t_1, \dots, t_n)$  correspond à une **décomposition** d'un arbre en un contexte et une suite d'arbres associés aux arguments du contexte.

Nous allons décomposer un arbre de sorte visible  $t \in T(\Omega)_{\Xi'}$  =  $\{t \in T(\Omega) \mid t^G \in \Xi'\}$  sous la forme  $t = C(t_1, \dots, t_n)$  où  $C$  est le plus petit contexte non trivial (i.e.,  $C \neq [\ ]$ ) pour lequel  $t_i^G \in \Xi'$ . En particulier si tous les sous arbres stricts  $t'$  de  $t$  sont tels que  $t' \notin \Xi'$ , on aura  $t = C$  (et donc le contexte ne contiendra aucun argument). Cette décomposition est donnée par la relation  $t \gg_C t_1 \dots t_n$  définie ci-dessous.

**Définition 3.5** (Décomposition d'un arbre). La relation  $\vdash \subseteq T(\Omega) \times \Omega^{\otimes} \times T(\Omega)^*$ , où  $(t, C, t_1 \dots, t_n) \in \vdash$  est noté  $t \vdash_C t_1 \dots t_n$ , est la plus petite relation telle que

- (1)  $t \vdash_{[\ ]} t$  si  $t^G \in \Xi'$
- (2)  $t \vdash_{\omega_n} t_1 \dots t_n$  si  $t = \omega(t_1, \dots, t_n)$  et  $t^G \notin \Xi'$
- (3) Si  $t \vdash_{\omega_n} t_1 \dots t_n$  et  $t_i \vdash_{C_i} v_i$  pour tout  $1 \leq i \leq n$  alors  $t \vdash_{\omega_n(C_1, \dots, C_n)} v$  où  $v = v_1 \cdots v_n$  est la concaténation des  $v_i$ .

Pour  $t = \omega(t_1, \dots, t_n) \in T(\Omega)_{\Xi'}$  on pose  $t \gg_{\omega_n(C_1, \dots, C_n)} v$  lorsque  $t_i \vdash_{C_i} v_i$  où  $v_i \in (T(\Omega)_{\Xi'})^*$  et  $v = v_1 \cdots v_n$  est la concaténation des suites d'arbres  $v_1, \dots, v_n$ . On pose  $t \gg v$  si  $t \gg_C v$  pour un certain contexte  $C$ .

La correspondance entre cette décomposition d'un arbre et la relation de dérivation sur les sortes est donnée par les lemmes suivants.

**Lemme 3.6.**  $t \vdash_C t_1 \cdots t_n \implies t^G \vdash^* t_1^G \cdots t_n^G$

*Démonstration.* La démonstration s'effectue par récurrence sur la preuve de la relation  $t \vdash_C t_1 \cdots t_n$ . Le cas de base se décline suivant la valeur de  $t^G$  :

$t^G \in \Xi'$ . On a d'une part  $t \vdash_{[\ ]} t$  et d'autre part  $t^G \vdash^* t^G$  par réflexivité de cette relation.

$t^G \notin \Xi'$ . Si  $t$  est de la forme  $t = \omega(t_1, \dots, t_n)$  on déduit que  $t \vdash_{\omega_n} t_1 \dots t_n$ . D'autre part  $t_1^G \dots t_n^G \in \mathcal{L}_{\omega, t^G}$  et donc  $t^G \vdash^* t_1^G \dots t_n^G$ .

Pour le cas général supposons  $t = \omega(t_1, \dots, t_n)$  et  $t_i \vdash_{C_i} v_i$  pour tout  $1 \leq i \leq n$ . On a dans ce cas  $t \vdash_{\omega_n(C_1, \dots, C_n)} v_1 \dots v_n$ . Par hypothèse de récurrence  $t_i^G \vdash^* v_i^G$  en convenant que  $(t'_1 \dots t'_k)^G = (t'_1)^G \dots (t'_k)^G$ . Par ailleurs de  $t_1^G \dots t_n^G \in \mathcal{L}_{\omega, t^G}$  il vient  $t^G \vdash^* t_1^G \dots t_n^G$ , et donc  $t^G \vdash^* v_1^G \dots v_n^G$ .  $\square$

Le lemme suivant en est une réciproque partielle où on se restreint aux dérivations maximales, c'est-à-dire que les sortes des arbres  $t_1, \dots, t_n$  sont visibles (la sorte  $s$  par contre ne le sera généralement pas).

**Lemme 3.7.**  $s \vdash_{\max}^* t_1^G \cdots t_n^G \implies \exists t \in T(\Omega) \exists C \in \Omega_n^{\textcircled{a}} \quad t \vdash_C t_1 \cdots t_n \text{ et } t^G = s$

*Démonstration.* La preuve est par récurrence sur la longueur de la dérivation de  $\vdash^*$ .

**Dérivations de longueur nulle :** ceci correspond au cas où  $s \in \Xi'$  —puisque la dérivation est maximale— avec  $n = 1$  et  $s = t_1^G$ . Le résultat découle de  $t_1 \vdash_{[]} t_1$ .

**Cas général :** la dérivation  $s \vdash^* t_1^G \cdots t_k^G$  peut se décomposer en  $s \vdash s_1 \dots s_n$  —ce qui signifie que  $s_1 \dots s_n \in \mathcal{L}_{\omega, s}$  pour un certain opérateur  $\omega \in \Omega$ — et, pour tout  $1 \leq i \leq n$ ,  $s_i \vdash v_i$  tels que  $v_1 \cdots v_n = t_1^G \cdots t_k^G$ . Puisque les  $t_i^G$  sont visibles ces dérivations sont maximales. Par hypothèse de récurrence il existe un arbre  $t_i \in T(\Omega)$  et un contexte  $C_i \in \Omega^{\textcircled{a}}$  tels que  $t_i^G = s_i$  et  $t_i \vdash_{C_i} v_i$ . L'arbre  $t = \omega(t_1, \dots, v_n)$  vérifie  $t^G = s$  et  $t \vdash_{\omega_n(C_1, \dots, C_n)} t_1 \cdots t_n$ .  $\square$

**Lemme 3.8.**  $(\exists s_1 \dots s_n \in \mathcal{L}_{\omega, s} \text{ t.q. } s \in \Xi' \text{ et } s_1 \dots s_n \vdash_{\max}^* (t'_1)^G \dots (t'_k)^G) \implies (\exists t \in T(\Omega) \text{ t.q. } t \gg t'_1 \dots t'_k \text{ et } t^G = s)$

*Démonstration.* Supposons  $s_1 \dots s_n \vdash_{\max}^* (t'_1)^G \dots (t'_k)^G$  pour  $s \in \Xi'$  et  $s_1 \dots s_n \in \mathcal{L}_{\omega, s}$ . Par le lemme 3.7 on déduit l'existence d'arbres  $t_i$  et de contextes  $C_i$  tels que  $t_i \vdash_{C_i} v_i$ ,  $v_1 \dots v_n = t'_1 \dots t'_k$  et  $t_i^G = s_i$ . Posons  $t = \omega(t_1, \dots, t_n)$ . De  $t_i \vdash_{C_i} v_i$ ,  $v_1 \dots v_n = t'_1 \dots t'_k$  on déduit  $t \gg t'_1 \dots t'_k$ . Par ailleurs de  $s_1 \dots s_n \in \mathcal{L}_{\omega, s}$  et  $t_i^G = s_i$  on déduit  $t^G = s$ .  $\square$

Tout arbre  $t \in T(\Omega)_{\Xi'}$  dont la sorte est visible peut se décomposer sous la forme  $t = C(t_1, \dots, t_n)$  dans lequel  $C$  est un contexte non trivial (i.e.,  $C \neq []$ ) et  $t_1, \dots, t_n$  sont des arbres de sortes visibles. La relation  $t \gg_C t_1 \cdots t_n$  correspond à une telle décomposition, plus précisément c'est celle pour laquelle le contexte  $C$  est minimal. Le contexte  $C$  et les arbres  $t_1, \dots, t_n$  tels que  $t \gg_C t_1 \cdots t_n$  sont donc déterminés de façon unique à partir de  $t$ . Autrement dit la relation  $\gg$  détermine une application  $\gg: T(\Omega)_{\Xi'} \rightarrow (T(\Omega)_{\Xi'})^*$ . Lorsque  $t \gg t_1 \cdots t_n$  les arbres  $t_1, \dots, t_n$  sont strictement plus petits que  $t$  —où la taille d'un arbre est donnée par son nombre de noeuds— et donc si  $P \subseteq T(\Omega)_{\Xi'}$ , vue comme une propriété des arbres dont la sorte est visible, vérifie les conditions suivantes <sup>1</sup>

- 1)  $t \gg \varepsilon \implies P(t)$
- 2)  $t \gg t_1 \cdots t_n \wedge P(t_1) \wedge \dots \wedge P(t_n) \implies P(t)$

alors  $P = T(\Omega)_{\Xi'}$ . Ceci nous donne un principe d'induction qui nous permet de raisonner par récurrence sur l'ensemble des arbres dont la sorte est visible en utilisant la décomposition associée à la relation  $\gg$ . On note par ailleurs que

$$t \gg_{\omega} t'_1 \dots t'_k \implies \pi(t) = \omega(\pi(t'_1) \cdots \pi(t'_k))$$

<sup>1</sup>On peut voir le premier cas comme un particulier du second si on admet le cas  $n = 0$  en convenant que la liste  $t_1 \cdots t_n$  représente alors la liste vide. La même convention s'applique pour la définition inductive des arbres. Avec cette convention nous évitons de considérer le cas de base dans la récurrence sur la décomposition par  $\gg$  ainsi que dans la récurrence sur la structure des arbres puisque à chaque fois les cas de bases apparaissent comme cas particuliers du cas général.

Nous disposons maintenant de tous les ingrédients nous permettant d'établir le théorème 3.3.

*Démonstration. [du théorème 3.3]*

Nous montrons que  $L(G/\pi, s) = \pi(L(G, s))$  pour tout  $s \in \Xi'$ . C'est-à-dire que les deux ensembles suivants coïncident  $\{\pi(t) \mid t^G = s\} = \{t \mid s \in t^{G/\pi}\}$ . Pour cela nous montrons successivement que (i)  $t^G \in (\pi(t))^{G/\pi}$  lorsque  $t^G \in \Xi'$  et (ii)  $s \in t^{G/\pi} \implies (\exists t \text{ t.q. } t^G = s \text{ et } \pi(t) = t')$ .

- (i) On raisonne par récurrence sur la décomposition de  $t$  donnée par la relation  $\gg$ . Supposons  $t \gg t'_1 \dots t'_k$ . Cela signifie que  $t = \omega(t_1, \dots, t_n)$  tel que  $t_1^G \dots t_n^G \in \mathcal{L}_{\omega, t^G}$ ,  $t_i \vdash_{C_i} v_i$  avec  $v_i \in (T(\Omega)_{\Xi'})^*$  pour  $1 \leq i \leq n$  et  $v_1 \dots v_n = t'_1 \dots t'_k$ . Par le lemme 3.6 il vient que  $t_i^G \vdash^* v_i^G$  pour  $1 \leq i \leq n$ . Par hypothèse de récurrence on déduit que  $(t'_i)^G \in \pi(t_i)^{G/\pi}$  et donc  $(t'_1)^G \dots (t'_k)^G \in \mathcal{A}^d(\omega, t^G)$ . Du fait que  $\pi(t) = \omega(\pi(t_1), \dots, \pi(t_k))$  il vient  $t^G \in \pi(t)^{G/\pi}$ .
- (ii) On raisonne par récurrence sur la structure de l'arbre  $t'$ . Supposons  $t' = \omega(t'_1 \dots t'_k)$  et  $s \in (t')^{G/\pi}$ . En utilisant la définition 3.2 on en déduit l'existence de  $s_1 \dots s_n \in \mathcal{L}_{\omega, s}$  et de  $s'_i \in (t'_i)^{G/\pi}$  pour tout  $1 \leq i \leq k$  tels que  $s_1 \dots s_n \vdash_{\max}^* s'_1 \dots s'_k$ . Par hypothèse de récurrence il existe des arbres  $t_i$  tels que  $t_i^G = s'_i$  et  $\pi(t_i) = t'_i$ . De  $s_1 \dots s_n \in \mathcal{L}_{\omega, s}$  et  $s_1 \dots s_n \vdash_{\max}^* t_1^G \dots t_k^G$  on en déduit par le lemme 3.8 l'existence d'un arbre  $t \in T(\Omega)_{\Xi'}$  tel que  $t \gg t_1 \dots t_k$  et  $t^G = s$ . De  $t \gg t_1 \dots t_k$  et  $\pi(t_i) = t'_i$  on déduit  $\pi(t) = \omega(\pi(t_1), \dots, \pi(t_k)) = \omega(t'_1 \dots t'_k) = t'$ .

Comme  $\# \in \Xi'$  le résultat découle de

$$\begin{aligned} \pi(L(G)) &= \{\pi(t_1) \dots \pi(t_n) \mid \#(t_1, \dots, t_n) \in L(G, \mathbf{ax})\} \\ &= \{\pi(t_1) \dots \pi(t_n) \mid \#(\pi(t_1), \dots, \pi(t_n)) \in L(G/\pi, \mathbf{ax})\} \\ &= L(G/\pi) \end{aligned}$$

Car  $\pi : L(G, s) \rightarrow L(G/\pi, s)$  est une surjection pour tout  $s \in \Xi'$ . □

## 4. Discussion

Nous nous intéressons aux deux problèmes suivants.

**PROBLÈME. — Invocation d'un service** – Supposons qu'on veuille développer un système, manipulant des documents conformes à une grammaire  $G = (\Omega, \Xi, \mathcal{L})$ , en utilisant des services dont l'interface est donnée par une autre grammaire  $G'$  qui décrit les documents que le service est disposé à recevoir. Un sous ensemble  $\Xi' \subseteq \Xi$  de sortes visibles peut représenter les parties du document qu'on souhaite transmettre au service pour traitement. On doit pouvoir vérifier que  $p_{\Xi'}(L(G)) \subseteq L(G')$  pour s'assurer que les documents transmis par le système au service seront toujours reconnus conformes par ce dernier.

**PROBLÈME. — Opacité** – Si un ensemble reconnaissable d'arbres  $S \subseteq T(\Omega)$  représente un secret vis-à-vis de l'observateur associé à la fonction d'abstraction  $\phi = p_{\Xi'}$  —opérant sur les arbres conformes à une grammaire  $G = (\Omega, \Xi, \mathcal{L})$  avec  $\Xi' \subseteq \Xi$ — alors  $S$  est

opaque vis-à-vis de  $\phi$  si et seulement si  $p_{\Xi'}(S \cap L(G)) \leq p_{\Xi'}(L(G) \setminus S)$ .

Le langage  $\mathcal{A}_{\omega,s} = \{v \in \Xi'^* \mid \exists u \in \mathcal{L}_{\omega,s} \quad u \vdash^* v\}$  associé à l'algèbre  $\mathcal{A} = (G/p_{\Xi'})$  qui caractérise les projections des documents conformes à la grammaire  $G$ , i.e.,  $L(\mathcal{A}) = p_{\Xi'}(L(G))$ , est un *langage algébrique*. Ces langages sont rationnels si la grammaire est *non récursive*, c'est-à-dire si on ne peut trouver de cycles dans la relation de dépendance  $s > s'$  correspondant au fait que  $s'$  apparaît en partie droite d'une production associée à  $s$ . Dans ce cas l'algèbre  $\mathcal{A}$  est une grammaire. Les grammaires sont closes de manière effective par les opérations booléennes. La complémentation et donc l'union ne préservent pas la non-circularité. Néanmoins les deux opérations qui nous intéressent, à savoir l'intersection et la différence, préservent cette propriété. Enfin on sait décider de la vacuité et donc de l'inclusion des langages de telles grammaires. Les deux problèmes ci-dessus peuvent donc être résolus de manière effective dans ce cas. En particulier on peut construire la grammaire  $\pi_{\Xi'}(S \cap G) \setminus p_{\Xi'}(G \setminus S)$  qui reconnaît exactement les cas de violation du secret. Dans le cadre d'un outil interactif cette grammaire peut être utilisée pour fournir des contre-exemples —les témoins— lorsque l'opacité n'est pas assurée. Cette information peut être utilisée par le concepteur du système pour l'aider à modifier la description du workflow en vue d'en assurer l'opacité.

Il est assez naturel de considérer comme nous l'avons fait des grammaires dans lesquelles les parties droites des productions sont des langages réguliers (et non un ensemble fini de mots) car on aura souvent besoin de considérer qu'un noeud d'un artefact puisse contenir une liste de sous documents d'une certaine sorte sans que la quantité de tels documents puissent fixée *a priori*. En revanche on aura rarement besoin de considérer des grammaires récursives. Ainsi du point de vue des artefacts l'hypothèse de non circularité est tout à fait raisonnable. Néanmoins comme nous l'avons signalé plus haut dans le texte cette hypothèse restreindra la classe des secrets que nous pourrions prendre en compte puisqu'on ne pourra pas exprimer des propriétés qui dépendent d'information contextuelles.

La solution alternative que nous considérons dans la section suivante est d'interpréter les documents modulo permutations des sous arbres d'un noeud. Nous restons de cette façon dans le cadre rationnel grâce au théorème de Parikh qui nous dit que l'image commutative d'un langage algébrique est rationnel.

---

## 5. Vérifier l'opacité d'un système

Nous considérons que l'ordre dans lequel les successeurs d'un noeud apparaissent n'est pas significatif. Pour prendre en compte cette contrainte on pourrait décider de se restreindre aux grammaires pour lesquelles les langages  $\mathcal{L}_{\omega,A}$  sont clos par commutations (la plus petite congruence  $\approx$  pour laquelle  $\omega \cdot \omega' \approx \omega' \cdot \omega$ ). Mais ceci peut être beaucoup trop restrictif. Pour cette raison, nous allons plutôt considérer la relation de "conformité modulo permutations" définie comme suit.

**Définition 5.1** (Conformité modulo permutations). *La relation de conformité modulo permutations est la plus petite relation telle que :*

$$(\forall i \in \{1, \dots, n\} \quad G \vdash_c t_i :: s_i \quad \wedge \quad s_1 \cdots s_n \in [\mathcal{L}_{\omega,s}]) \Rightarrow G \vdash_c \omega(t_1, \dots, t_n) :: s$$

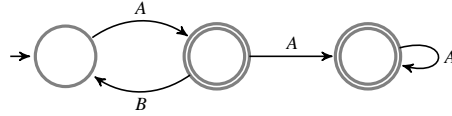
où  $[L] = \{u \mid \exists v \in L \ u \approx v\}$  est la clôture de  $L$  modulo commutations. Le langage modulo commutations d'une grammaire est donné par :

$$\begin{aligned} L_c(G, s) &= \{t \in T(\Omega) \mid G \vdash_c t :: s\} \\ L_c(G) &= \{t_1 \dots t_n \in T(\Omega) \mid G \vdash_c \#(t_1, \dots, t_n) :: \mathbf{ax}\} \end{aligned}$$

REMARQUE. — Pour que la procédure de reconnaissance soit déterministe nous supposons, non seulement que les différents langages  $\mathcal{L}_{\omega, s}$  —lorsque  $s$  varie— sont disjoints, mais que leurs clôtures commutatives le sont également :  $s \neq s' \Rightarrow [\mathcal{L}_{\omega, s}] \cap [\mathcal{L}_{\omega, s'}] = \emptyset$ .

Dans la mesure où la clôture par commutations d'un langage rationnel n'est généralement pas un langage rationnel, cette seconde approche est plus générale.

**Exemple 5.2.** La clôture par permutation de  $\mathcal{L} = (AB)^* \cdot A^*$  est l'ensemble des mots dont le nombre d'occurrences de  $A$  est supérieur ou égal au nombre d'occurrences de  $B$ . Ce dernier n'est pas rationnel (il possède un ensemble infini de résidus). On peut utiliser l'automate associé au langage régulier  $\mathcal{L}$



comme un générateur de  $[\mathcal{L}]$ . Pour cela on interprète l'expression rationnelle  $(AB)^* A^*$  dans  $\mathbb{N}^2$  en identifiant les lettres  $A$  et  $B$  avec les deux générateurs  $(1, 0)$  et  $(0, 1)$  respectivement. Au niveau de l'automate cela revient à cumuler les valeurs des étiquettes rencontrées le long d'un chemin reconnaissant de l'automate. Cet automate procure ainsi un mécanisme qui permet d'engendrer  $[\mathcal{L}] = \{(x, y) \in \mathbb{N}^2 \mid x \geq y\}$ . Cet ensemble est associé à la formule  $F(x, y) \equiv (\exists z) \cdot x = y + z$ . Néanmoins cet automate ne peut être utilisé comme un reconnaiseur.

RAPPEL. — On peut définir, de façon générale, ce qu'est une partie reconnaissable d'un monoïde. Les parties rationnelles et reconnaissables coïncident dans le cas des monoïdes libres mais pas en général. Les parties reconnaissables d'un monoïde finiment engendré sont rationnelles. A titre d'exemple les parties reconnaissables de  $\mathbb{N}^2$  sont les unions finies de "grilles rectangulaires". Ces dernières désignent des ensembles de la forme  $\{(x_0 + n \cdot x_1, y_0 + m \cdot y_1) \mid n, m \in \mathbb{N}\}$  où  $x_0, x_1, y_0$  et  $y_1$  sont des constantes entières et on peut vérifier que le quadrant  $\{(x, y) \in \mathbb{N}^2 \mid x \geq y\}$  ne peut s'écrire comme une union finie de telles grilles (il s'agit d'une partie rationnelle non reconnaissable).

Les ensembles qui nous intéressent sont donc les images commutatives des langages rationnels, c'est-à-dire les parties rationnelles de  $\mathbb{N}^k$  (où  $k$  est la taille de l'alphabet). Le fait que ces ensembles ne soient pas reconnaissables pourrait sembler constituer un obstacle pour la vérification de la conformité d'un document (modulo commutations). En fait il n'en est rien, bien que ce processus va devoir suivre une procédure un peu plus compliquée que dans le cas de base. Les parties rationnelles de  $\mathbb{N}^k$  coïncident avec les parties



semi-linéaires.

RAPPEL. — Une partie linéaire est l'ensemble des vecteurs de la forme  $v_0 + n_1 \cdot v_1 + \dots + n_\ell \cdot v_\ell$  (une expression affine dans laquelle les vecteurs  $v_0$  —l'origine— et  $v_1, \dots, v_\ell$  —les vecteurs de base— sont fixés et les variables  $n_i$  représentent des entiers arbitraires). Une partie semi-linéaire est une union finie de parties linéaires.

Ces ensembles semi-linéaires coïncident par ailleurs avec les ensembles définissables par des formules de Presburger, i.e., de la forme  $\{(x_1, \dots, x_n) \mid F(x_1, \dots, x_n)\}$  où  $F$  est une formule de Presburger.

RAPPEL. — La logique de Presburger est le fragment de l'arithmétique dans lequel on exclut la multiplication, c'est-à-dire qu'il s'agit du calcul propositionnel avec quantification sur les entiers et avec l'addition. On peut bien sûr écrire la multiplication par une constante puisque par ex.  $3x = x + x + x$ . De la même manière on peut écrire des inéquations : par ex.  $x \leq y$  est une abréviation de  $(\exists z) \cdot y = x + z$ .

Ces différentes correspondances sont effectives : on peut passer entre expressions régulières, ensembles semi-linéaires et formules de Presburger en utilisant des algorithmes bien établis (voir [10, 11, 14, 6]). L'intérêt de la logique de Presburger est sa décidabilité (au contraire de l'arithmétique dans son ensemble, qui est indécidable). Comme pour la décision de la logique du premier ordre, la méthode consiste à construire (par induction sur la structure de la formule) un automate qui reconnaît l'ensemble des vecteurs vérifiant la formule (pour un codage particulier des vecteurs en des mots sur un certain alphabet). La validité de la formule se réduit en la non vacuité de l'automate correspondant. Nous utilisons cette construction pour vérifier comme suit la validité d'un arbre. Considérant un noeud étiqueté par un opérateur  $\omega$ , on suppose inductivement que chacun de ses sous-arbres immédiats ait été jugé conforme avec une sorte déterminée, on collecte ces différentes sortes pour former un vecteur dont on vérifie qu'il satisfait la formule de Presburger associée à l'image commutative de  $\mathcal{L}_\omega = \cup_{s \in \Xi} \mathcal{L}_{\omega, s}$ .

Les ensembles semi-linéaires sont par ailleurs clos de manière effective par les opérations booléennes, ce qui fait que, de façon analogue à ce que nous avons exprimé plus haut, nous pouvons décider si le langage d'une grammaire est à commutations près contenu dans celui d'une autre, i.e., la relation  $G \leq_c G' \iff L_c(G) \subseteq L_c(G')$  est décidable.

La définition suivante est l'adaptation de la définition 3.2 pour les grammaires à commutations près.

**Définition 5.3.** Si  $\Xi' \subseteq \Xi$  est un sous-ensemble de l'alphabet des sortes d'une grammaire  $G = (\Omega, \Xi, \mathcal{L})$  et  $s \in \Xi \setminus \Xi'$ , on pose  $L(G, \Xi', s) = \{u \in \Xi'^* \mid s \vdash^* u\}$ . La projection de  $G$  à  $\Xi'$  est la grammaire  $p_{\Xi'}(G) = (\Omega, \Xi', \mathcal{L}')$  où  $\mathcal{L}'_{\omega, s'} = \mathcal{L}_{\omega, s'}[[L(G, \Xi', s)]]/s$  ;  $s \in \Xi \setminus \Xi'$ , c'est-à-dire qu'on substitue l'image commutative du langage algébrique  $L(G, \Xi', s)$  (qui est donc un langage rationnel par le théorème de Parikh) à la variable non visible  $s \in \Xi \setminus \Xi'$  dans la partie droite  $\mathcal{L}_{\omega, s'}$  de la règle, dans la grammaire d'origine, associée au symbole visible  $s' \in \Xi'$ .

REMARQUE. — Dans la pratique nous n'introduisons pas explicitement le langage régulier  $[L(G, \Xi', s)]$  mais on exhibe une expression régulière qui le caractérise à commutations près. Pour cela on utilise le fait que la plus petite solution d'une équation  $X = A(X) \cdot X + T$  dans laquelle la variable  $X$  n'apparaît pas dans  $T$  est à commu-

tations près donnée par l'expression régulière  $A(T)^* \cdot T$

**Exemple 5.4** (voir figure 2).

*La projection de la grammaire*

$$G = \left( \begin{array}{l} A \text{ where} \\ A \rightarrow \omega_1\langle BC \rangle + \omega_2\langle \varepsilon \rangle \\ B \rightarrow \omega_3\langle BD \rangle + \omega_4\langle \varepsilon \rangle \\ C \rightarrow \omega_1\langle CA \rangle + \omega_4\langle CCB \rangle + \omega_1\langle D \rangle \\ D \rightarrow \omega_2\langle AB \rangle \end{array} \right)$$

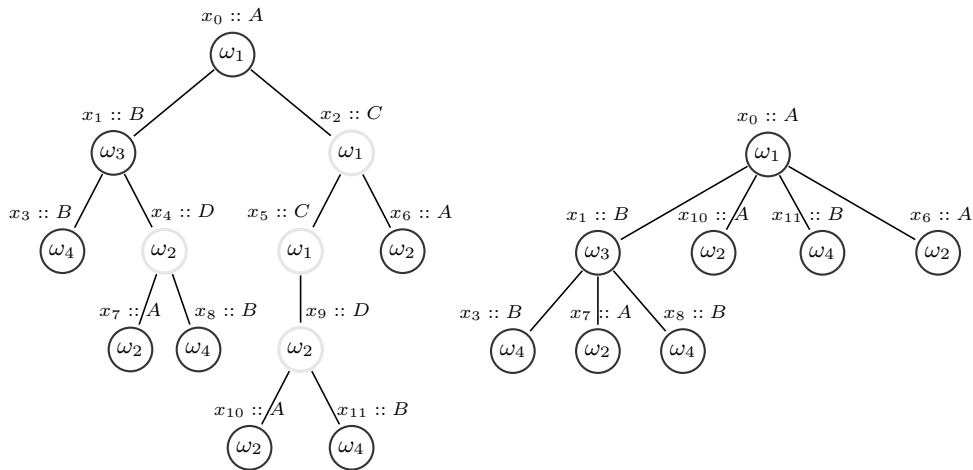
sur le sous-alphabet  $\Xi' = \{A, B\}$  est donnée par

$$p_{\{A,B\}}(G) = \left( \begin{array}{l} A \text{ where} \\ A \rightarrow \omega_1\langle BC \rangle + \omega_2\langle \varepsilon \rangle \\ B \rightarrow \omega_3\langle BD \rangle + \omega_4\langle \varepsilon \rangle \\ C =_c CA + CCB + D \\ D =_c AB \end{array} \right)$$

En utilisant  $CA + CCB + D =_c (A + CB)C + D$  on obtient  $C =_c (A + DB)^*D = (A + ABB)^*AB$  et donc

$$p_{\{A,B\}}(G) = \left( \begin{array}{l} A \text{ where} \\ A \rightarrow \omega_1\langle B(A + ABB)^*AB \rangle + \omega_2\langle \varepsilon \rangle \\ B \rightarrow \omega_3\langle BAB \rangle + \omega_4\langle \varepsilon \rangle \end{array} \right)$$

L'arbre de la figure 1 (qu'on retrouve à gauche de la figure 2) se présente sous la forme



**Figure 2.** A gauche, l'arbre  $t$  conforme à la grammaire  $G$  provenant de la figure 2 :  $G \vdash_c t :: A$ . A droite, la projection de  $t$  sur le sous-alphabet de sortes visibles  $\Xi' = \{A, B\}$ . On vérifie que  $p_{\{A,B\}}(G) \vdash_c p_{\{A,B\}}(t) :: A$ .

du système d'équations suivant :

$$\begin{array}{llll} x_0 = \omega_1(x_1, x_2) & x_3 = \omega_4() & x_6 = \omega_2() & x_9 = \omega_2(x_{10}, x_{11}) \\ x_1 = \omega_3(x_3, x_4) & x_4 = \omega_2(x_7, x_8) & x_7 = \omega_2() & x_{10} = \omega_2() \\ x_2 = \omega_1(x_5, x_6) & x_5 = \omega_1(x_9) & x_8 = \omega_4() & x_{11} = \omega_4() \end{array}$$

On peut étiqueter chaque variable par la sorte qui lui est associée lorsqu'on résout ce système dans l'algèbre associée à la grammaire  $G$

$$\begin{array}{llll} x_0^A = \omega_1(x_1^B, x_2^C) & x_3^B = \omega_4() & x_6^A = \omega_2() & x_9^D = \omega_2(x_{10}^A, x_{11}^B) \\ x_1^B = \omega_3(x_3^B, x_4^D) & x_4^D = \omega_2(x_7^A, x_8^B) & x_7^A = \omega_2() & x_{10}^A = \omega_2() \\ x_2^C = \omega_1(x_5^C, x_6^A) & x_5^C = \omega_1(x_9^D) & x_8^B = \omega_4() & x_{11}^B = \omega_4() \end{array}$$

Le système d'équation associé à la projection de l'arbre est obtenu en "effaçant" le symbole d'opérateur en partie droite d'une équation dont la variable en partie gauche est d'une sorte non visible :

$$\begin{array}{llll} x_0^A = \omega_1(x_1^B, x_2^C) & x_3^B = \omega_4() & x_6^A = \omega_2() & x_9^D = x_{10}^A, x_{11}^B \\ x_1^B = \omega_3(x_3^B, x_4^D) & x_4^D = x_7^A, x_8^B & x_7^A = \omega_2() & x_{10}^A = \omega_2() \\ x_2^C = x_5^C, x_6^A & x_5^C = x_9^D & x_8^B = \omega_4() & x_{11}^B = \omega_4() \end{array}$$

On simplifie ensuite ce système en supprimant toute équation de la forme  $x^s = x_1^{s_1} \dots x_n^{s_n}$  associée à une sorte non visible  $s \in \Xi \setminus \Xi'$  après avoir substitué la définition de  $x_s$ , i.e.,  $x_1^{s_1} \dots x_n^{s_n}$ , à chacune de ses occurrences d'utilisation (en fait ici on ne peut avoir qu'une seule telle occurrence car ce système d'équation provient d'un arbre) :

$$\begin{array}{llll} x_0^A = \omega_1(x_1^B, x_{10}^A, x_{11}^B, x_6^A) & x_3^B = \omega_4() & x_7^A = \omega_2() & x_{10}^A = \omega_2() \\ x_1^B = \omega_3(x_3^B, x_7^A, x_8^B) & x_6^A = \omega_2() & x_8^B = \omega_4() & x_{11}^B = \omega_4() \end{array}$$

## 6. conclusion

Dans ce papier nous avons proposé une formalisation de l'opacité des artefacts dans les systèmes workflow centrés sur l'échange de données structurées. Nous avons identifié deux contextes dans lesquels non seulement l'opacité est décidable mais l'algorithme de décision peut générer des contre exemples à l'opacité lorsque cette propriété n'est pas satisfaite : (i) lorsque les grammaires sont non récursives ou (ii) lorsque les documents ne sont considérés qu'à la permutation près de sous arbres.

Le retour d'information produit par l'algorithme de décision peut être utilisé par le concepteur du système pour opérer les modifications nécessaires au rétablissement de l'opacité. On peut souhaiter que le système aille un peu plus loin en effectuant lui-même ces modifications ou tout du moins en proposant à l'utilisateur des modifications de nature à résoudre les éventuels problèmes de fuite d'information.

De telles techniques, dite de contrôle, ont été étudiées dans le contexte des systèmes à événement discrets, cadre dans lequel l'opacité a été étudiée jusque là. Rappels brièvement ce qui a été fait pour assurer l'opacité dans les systèmes à événement discrets. Dans ce cas l'univers  $U$  est donné par un langage régulier dont les éléments (des mots) sont toutes les exécutions du système, c'est-à-dire les suites d'événements (actions internes ou interactions avec un utilisateur) qui peuvent advenir dans le système. Chaque observateur

ne voit que les actions avec lesquelles il interagit avec le système (un sous alphabet de l'ensemble  $E$  des événements). La fonction d'observation qui lui est associée  $\phi : U \rightarrow O$  est la projection qui efface dans toutes les exécutions toutes les occurrences d'événements non observables. Comme (i) l'ensemble des langages réguliers est clos de manière effective par les opérations booléennes, (ii) l'image d'un langage régulier par un morphisme de monoïdes est (de manière effective) un langage régulier, (iii) on peut décider de la vacuité d'un langage régulier (et donc aussi de l'inclusion de deux langages réguliers) on en déduit qu'on peut décider si  $\phi(S) \subseteq \phi(U \setminus S)$ , c'est-à-dire de l'opacité de  $S$ .

Le langage régulier  $L_{\phi,S} = \phi^{-1}(\phi(S) \setminus \phi(U \setminus S))$  donne toutes les exécutions qui témoignent de la non-opacité de  $P$  vis-à-vis du secret  $S$ . De façon plus précise le langage  $L = U \setminus (\bigcup \{\phi^{-1}(\phi(S) \setminus \phi(U \setminus S)) \mid S \text{ secret pour } \phi\})$  est le plus grand langage opaque contenu dans  $U$ . Si on fait l'hypothèse que toutes les actions sont contrôlables on peut rendre le système opaque en le synchronisant avec un contrôleur donné par le langage précédant. C'est-à-dire qu'on autorise le système à exécuter une séquence que lorsque celle-ci est de façon simultanée autorisée, c'est-à-dire exécutable, par le contrôleur. Le langage du système contrôlé est donc l'intersection du langage du système d'origine avec celui du contrôleur. Ce contrôleur est optimal en ce sens qu'il n'interdit que les comportements qui violent l'opacité du système. Qui plus est [1] ce contrôle est distribué, c'est-à-dire qu'on peut de façon équivalente associer à chaque utilisateur  $\phi$  un contrôleur qui restreint les interactions autorisées de cet utilisateur avec le système. Des résultats analogues existent également sous des hypothèses plus faibles, c'est-à-dire lorsqu'on ne suppose plus que toutes les actions puissent être contrôlées.

Nous disposons ici de la même caractérisation logique du contrôleur quitte à remplacer des automates de mots par des automates d'arbres (non ordonnés et à arité variables). Par ailleurs nous savons [2] que la synchronisation de vues partielles d'un même document peut s'obtenir via le produit des automates d'arbres associés à chacune de ces vues. Les techniques de synthèse et de distribution d'un contrôleur devraient pouvoir s'étendre sans difficultés majeures aux artefacts d'un workflow.

La solution obtenue par contrôle produit une grammaire qui génère exactement les documents de  $U$  pour lesquels aucun secret n'est dévoilé. Le résultat peut néanmoins être peu satisfaisant d'un point de vue pratique si cela force à exclure trop de documents, et en particulier des documents que le concepteur du système souhaiterait conserver dans son modèle. Une solution moins intrusive, qui s'inspire de techniques d'extensions conservatives de schémas XML, est à l'étude. Cette méthode consiste à n'invalider aucun document mais au contraire à étendre la grammaire : l'introduction de nouveaux documents licites constitue du "bruit" qui peut empêcher la détection de secrets.

Une suite de ce travail va s'attacher à préciser la forme des règles sémantiques qui régissent le calcul des attributs aux noeuds des documents. Cela déterminera par la même occasion la syntaxe pour les algèbres et donc en particulier pour les requêtes faites sur un document et pour l'expression des secrets. Nous comptons également implémenter les algorithmes de décision et de contrôle de l'opacité en utilisant les bibliothèques existantes de manipulation des ensembles semi-linéaires.

## 7. Bibliographie

- [1] ERIC BADOUEL, MAREK A. BEDNARCZYK , ANDRZEJ M. BORZYSZKOWSKI, BENOÎT CAILLAUD , PHILIPPE DARONDEAU, « Concurrent Secrets », *Discrete Event Dynamic Systems*, vol. 17, n° 4, 2007 :425–446.
- [2] ERIC BADOUEL, MAURICE TCHOUPÉ TCHENDJI « Merging Hierarchically-Structured Documents in Workflow Systems », *Electronic Notes in Theoretical Computer Science*, vol. 203, n° 5, 2008 :3–24.
- [3] ERIC BADOUEL, RODRIGUE TCHOUGONG , CÉLESTIN NKUIMI-JUGNIA, BERNARD FOT-SING « Attribute grammars as tree transducers over cyclic representations of infinite trees and their descriptive composition », *Theoretical Computer Science*, vol. 480, n° 0, 2013 :1–25.
- [4] ANNE BRUGGEMANN-KLEIN, A. MAKOTO MURATA, DERICK WOOD, « Regular Tree and Regular Hedge Languages over Unranked Alphabets », *HKUST-TCSC-2001-05*, 2001.
- [5] JEREMY BRYANS, MACIEJ KOUTNY, LAURENT MAZARÉ, PETER Y. A. RYAN, « Opacity generalised to transition systems », *Int. J. Inf. Sec.*, vol. 7, n° 6, 2008 :421–435.
- [6] FABRICE CHEVALIER, JÉRÉMIE CHALOPIN, « Représentation et algorithmique des ensembles semi-linéaires », *Rapport de stage, LSV - ENS Cachan*, 2001.
- [7] DAVID COHN, RICHARD HULL, « Business Artifacts : A Data-centric Approach to Modeling Business Operations and Processes », *Bulletin of the IEEE Computer Society Technical Committee on Data Engineering*, 2009.
- [8] HUBERT COMON, MAX DAUCHET, RÉMI GILLERON, FLORENT JACQUEMARD, DENIS LUGIEZ, CHRISTOF LÖDING, SOPHIE TISON, MARC TOMMASI, « Tree Automata Techniques and Applications », <http://tata.gforge.inria.fr>, 2008.
- [9] JÉRÉMY DUBREIL, THIERRY JÉRON, HERVÉ MARCHAND, « Monitoring information flow by diagnosis techniques », *European Control Conference, ECC'09*, 2009.
- [10] SEYMOUR GINSBURGH, « The Mathematical Theory of Context-Free Languages », *Mc Graw-Hill*, 1966.
- [11] SEYMOUR GINSBURGH, E.H. SPANIER, « Semigroups, Presburger formulas, and languages », *Pacific Journal of Mathematics*, vol. 16, n° 2, 1966 :285–296. .
- [12] DOMINIC HUGHES, VITALY SHMATIKOV, « Information Hiding, Anonymity and Privacy : a Modular Approach », *Journal of Computer Security*, vol. 12, n° 1, 2004 :3–36.
- [13] FENG LIN, « Opacity of discrete event systems and its applications », *Automatica*, vol. 47, n° 3, 2011 :496–500.
- [14] CHRISTOPHE REUTENAUER, « Aspects mathématiques des réseaux de Petri », *Masson, Paris*, 1988.