

Contrôle des observations pour la gestion des systèmes de flux de données.

Christophe Dousson*, Pierre Le Maigat*

*Orange Labs – 2, avenue Pierre Marzin – 22300 Lannion
{christophe.dousson, pierre.lemaigat}@orange-ftgroup.com
<http://perso.rd.francetelecom.fr/dousson>

Résumé. Les systèmes d'analyse de flux de données prennent de plus en plus d'importance dans un contexte où les données circulant sur les réseaux sont de plus en plus volumineuses et où la volonté de réagir au plus vite, en temps réel, devient un besoin nécessaire. Afin de permettre des analyses aussi rapides et efficaces que possible, il convient de pouvoir contrôler les flots de données et de focaliser les traitements sur les données pertinentes. Le protocole présenté dans ce papier donne au module de traitement des capacités d'action et de contrôle sur les observations remontantes en fonction de l'état de l'analyse. La diminution des flux résultant de telles focalisations permet des traitements beaucoup plus efficaces, plus pertinents et moins consommateurs de ressources. Les premiers résultats montrent un réel gain de performances sur nos applications (facteur 100).

Nous proposons donc ici un protocole permettant de propager des informations de contrôle du plus haut-niveau de l'analyse jusqu'aux sources d'événements. L'architecture mise en œuvre, baptisée TESS (pour *Timestamped Event Stream System*) est de type « workflow » où les événements transitent de module en module par des « liens ». Ces liens sont orientés d'une interface dite « Producteur » vers une interface dite « Consommateur » (voir figure 1) sur lesquels vont circuler les données du flot. Cette architecture s'appuie sur les hypothèses suivantes :

- les événements sont tous instantanés (les informations avec durée pourront être modélisées avec un événement de début et un autre de fin),
- les événements sont tous datés (avec une date ponctuelle) et seront donc notés (e, t) ,
- les connexions entre producteur et consommateur sont de type FIFO (en revanche, il n'y a pas de contrainte sur le fonctionnement interne d'un module),
- les envois de messages et d'événements sont tous *asynchrones*.

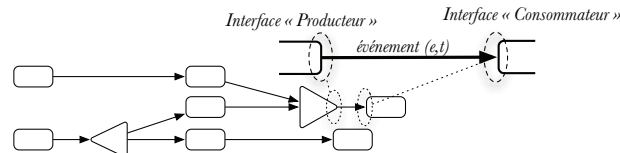


FIG. 1 – Architecture générale : producteurs et consommateurs

Protocole de contrôle

Un message de contrôle du flux d'événements est un *triplet* constitué *i*) d'un *type* qui correspond à l'action que l'on veut avoir sur le flux, *ii*) d'une *fenêtre temporelle* **tw** limitant l'action aux seuls événements dont la date d'occurrence est contenue dans celle-ci et *iii*) d'une *condition atemporelle* **C** limitant l'action aux seuls événements dont les données (atemporelles) vérifient cette condition (par exemple, une plage d'adresses IP).

Messages de contrôle du producteur vers le consommateur

NO_MORE_EVENT(**tw**, **C**). Ce message est un message de clôture du flux : *tous* les événements datés dans **tw** et vérifiant **C** ont été transmis. Il permet de gérer l'avancement du temps par émission régulière de `NO_MORE_EVENT ([-∞, t], TRUE)`.

MISSING_EVENT(**tw**, **C**). Permet de prévenir les consommateurs qu'il est possible que certains événements aient été perdus (volontairement ou non) par ou en amont du producteur.

RULE_ON(**tw**, **C**). Ce message est émis lorsqu'on souhaite statuer sur certains événements. En réponse à ce message, on s'attend à recevoir des *DISCARD* ou *FOCUS* (cf. ci-dessous).

Messages de contrôle du consommateur vers le producteur

FOCUS(**tw**, **C**). Tous les événements *présents ou à venir* vérifiant **tw** et **C** doivent être transmis dès que possible. Le producteur doit tout mettre en œuvre pour respecter cette urgence.

DISCARD(**tw**, **C**). Les événements *présents ou à venir* vérifiant **tw** et **C** n'ont plus d'utilité pour la suite de la chaîne de traitement ; ils peuvent donc être supprimés sans être transmis.

REMOVE_FOCUS(**tw**, **C**) (resp. **REMOVE_DISCARD**(**tw**, **C**)). Ces messages permettent de lever l'urgence (*FOCUS*) ou le filtre (*DISCARD*) précédemment postés.

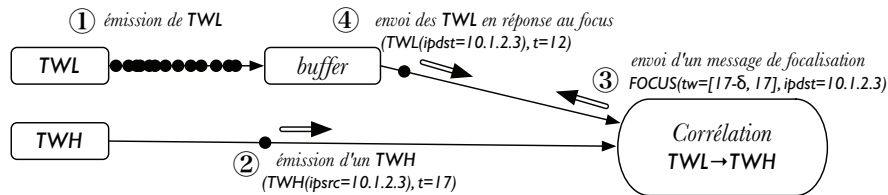


FIG. 2 – Exemple d'utilisation : corrélation avec focalisation.

Summary

Event Stream Processing (ESP) and Data Stream Management System (DSMS) become more and more popular in a world where huge amount of data flood the networks and where reactivity should as fast as possible. In order to be able to process efficient and quick on-line analysis, it is necessary to allow some control on the event streams and to focus the processing on relevant data. The protocol introduced in this paper is devoted to give control capabilities to the analysis module, depending on its current needs. The reduction of the flow rate induced by this allows more efficient and less resource consuming processing.