

Diffusion Vidéo avec une Meilleure Qualité d'Expérience et Respectant la Vie Privée

Simon Da Silva

1 Introduction

Avec l'essor d'Internet pour le grand public sont apparus de nouveaux modes de consommation. Tout d'abord, les contenus vidéo sur des plates-formes telles que YouTube, Vimeo ou Dailymotion augmentent en quantité et en qualité chaque année, notamment grâce à la rémunération des auteurs via la publicité et les partenariats. Il y a également une demande croissante du public pour des contenus spécialisés, dorénavant consommables sur tous supports (ordinateur, tablette, smartphone, téléviseur connecté, etc.). Par ailleurs, les vidéophiles souhaitent visionner des films et séries à leur rythme, sur leur support de prédilection, sans dépendre du programme télévisé. Pour cette raison, des services de vidéo à la demande (VoD) ont émergé. Ces méthodes de diffusion "Over The Top" grâce à des sites web ou des applications deviennent les solutions préférées des consommateurs et consommatrices.

En effet, le streaming vidéo représente actuellement plus de 60% du trafic Internet mondial total, 65% du trafic mobile sur Internet, et devrait atteindre 82% du trafic total d'ici 2022. Youtube et Netflix représentent à eux seuls plus de 50% du trafic Internet en heure de pointe aux États-Unis. De plus, l'essence même de la télévision (la diffusion en direct) est en train de basculer vers des alternatives sur Internet comme Twitch, privilégiées par le public grâce à leur flexibilité d'utilisation et de support. La part du streaming vidéo en direct sur Internet est également en très forte croissance, puisqu'elle devrait être multipliée par 15 pour atteindre 17% du trafic vidéo total d'ici 2022.

Ces nouveaux modes de consommation de contenus posent un problème majeur : la *qualité d'expérience* proposée à l'utilisateur. En effet, avec une telle croissance, les opérateurs et fournisseurs de contenu peinent à mettre à niveau leurs infrastructures pour supporter la demande toujours croissante des utilisateurs. Les réseaux sont souvent saturés, les serveurs se retrouvent surchargés, et il devient de plus en plus difficile de proposer une diffusion fiable, sans coupures, avec une bonne qualité visuelle et une stabilité satisfaisante, à des coûts abordables pour les fournisseurs de contenus. Il est alors nécessaire de trouver des solutions pour réduire l'impact de ces flux sur la santé du réseau, en permettant au plus grand nombre d'accéder aux ressources tout en fournissant une bonne qualité d'expérience aux utilisateurs consommant les contenus.

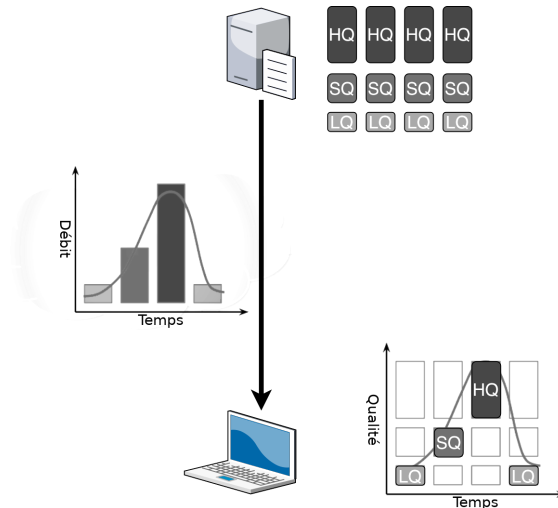


FIGURE 1 – Streaming adaptatif sur HTTP

Des méthodes de streaming vidéo adaptatif sur HTTP ont récemment vu le jour, telles que Adobe HDS, Apple HLS, Microsoft SS, et notamment le standard DASH qui est utilisé entre autres par YouTube, Facebook, Netflix et Twitch. L'objectif de ces techniques est de réduire le nombre de coupures lors de la diffusion de vidéos en adaptant la qualité du flux à la bande passante disponible entre l'utilisateur et le serveur (voir Figure 1). Pour cela, la vidéo est d'abord encodée dans plusieurs qualités. Ensuite, chaque qualité est découpée en segments de quelques secondes. Quand le client souhaite recevoir une vidéo, le serveur lui fournit une liste des différentes qualités disponibles, et le lecteur vidéo choisit alors la qualité la plus

adaptée à la bande passante disponible pour chaque segment. La vidéo est alors reçue en plusieurs segments qu'il faut remettre bout à bout pour lire le flux.

2 Motivation

Le standard DASH et les techniques similaires permettent d'améliorer sensiblement la *qualité d'expérience* du public en éliminant la plupart des coupures dues aux mauvaises conditions du réseau entre l'utilisateur et le serveur. En revanche, les problèmes liés à la surcharge des serveurs ou à leur capacité perdurent. Si de nombreux utilisateurs situés dans la même zone géographique regardent simultanément un même contenu vidéo, le serveur le plus proche devient rapidement surchargé (voir Figure 2). Certains utilisateurs subissent alors des dégradations de qualité ou une indisponibilité du contenu, et donc une *qualité d'expérience* faible et inéquitable.

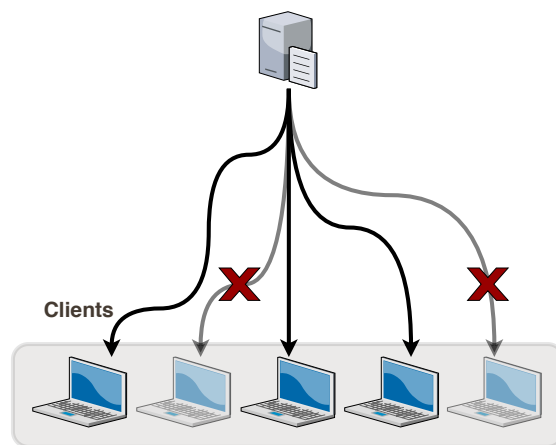


FIGURE 2 – Congestion d'un serveur DASH

Une autre problématique importante est la protection de la *vie privée*. L'utilisation des plateformes de streaming génère des informations personnelles sensibles en terme d'historique de visionnage aux vidéos. Ces données peuvent être exploitées soit au bénéfice de l'utilisateur, par exemple pour lui faire des recommandations personnalisées pour d'autres contenus, ou au bénéfice de la plateforme pour de la publicité ciblée. Cependant, la disponibilité des historiques d'accès peut également

conduire à des menaces majeures sur la vie privée. En effet, il est facilement possible d'inférer des informations privées sur l'utilisateur, tel que son genre, origine, ses orientations politiques, religieuses ou sexuelles, ou la composition du domicile familial.

Objectif

L'objectif de cette thèse est de proposer un système pragmatique et réaliste de streaming vidéo préservant la vie privée, proposant à la fois une meilleure *qualité d'expérience* et des garanties de protection de la *vie privée* aux utilisateurs, au moindre coût. Proposer une bonne *qualité d'expérience* signifie (1) fournir une qualité d'image haute, (2) minimiser les fluctuations de qualité, (3) éviter les interruptions pendant la lecture, et (4) assurer un temps de démarrage rapide. Protéger la *vie privée* des utilisateurs dans un système de streaming vidéo signifie camoufler leur historique de visionnage, à la fois des serveurs et des autres utilisateurs.

3 Contexte

Les solutions utilisant le standard DASH sont très efficaces pour faire face aux fluctuations de bande passante entre l'utilisateur et le serveur. En revanche, lorsque certains serveurs sont surchargés ou que des liens réseau sont saturés, il est judicieux de récupérer simultanément le contenu vidéo depuis plusieurs sources différentes afin de maximiser la qualité d'expérience de l'utilisateur. Plusieurs protocoles de streaming vidéo multi-sources ont récemment émergé pour faire face à la surcharge des serveurs ou liens réseau.

MS-STREAM

MS-STREAM (Multiple-Source Streaming), conçu au LaBRI, est un protocole de streaming vidéo adaptatif compatible avec DASH. Il permet d'utiliser plusieurs serveurs simultanément pour assurer une meilleure *qualité d'expérience* au public, à

la fois en réduisant le nombre de coupures et en améliorant la qualité vidéo affichée (grâce à l'agrégation des bandes passantes).

Tout comme pour DASH, la vidéo est d'abord encodée en différentes qualités puis découpée en segments contenant plusieurs groupes d'images. Le contenu est ensuite copié sur plusieurs serveurs différents. Lorsque l'utilisateur souhaite regarder un segment d'un contenu vidéo, le lecteur vidéo fait des requêtes auprès des différents serveurs disponibles). Chaque serveur va alors proposer un sous-segment composé de groupes d'images en bonne qualité et d'autres en qualité basse, en fonction de la bande passante disponible et de sa capacité. De cette manière, le client peut rassembler les différents groupes d'images reçus pour reformer un segment en bonne qualité. Si certains groupes d'images en bonne qualité ne sont pas reçus, il est possible d'utiliser ceux de basse qualité fournis par les autres serveurs afin de compléter le segment et continuer la lecture sans interruption.

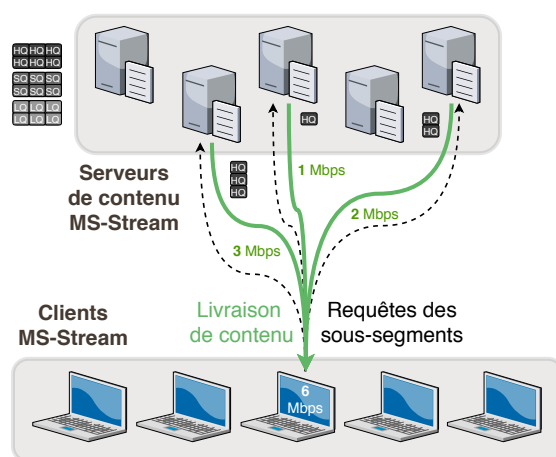


FIGURE 3 – Agrégation de bande passante avec MS-STREAM

Dans l'exemple de la Figure 3, l'utilisateur a une bande passante de 3 Mbps avec un serveur, 1 Mbps avec le deuxième, et 2 Mbps avec le troisième, une qualité visuelle allant jusqu'à 6 Mbps pourra être obtenue. Cette qualité est alors supérieure à la qualité que DASH aurait pu fournir (ici au maximum 3 Mbps avec le premier serveur). Si maintenant le deuxième serveur devient surchargé ou indisponible, l'utilisateur pourra toujours obtenir $3+2 = 5$ Mbps depuis les deux autres serveurs, sans que cela n'interrompe la lecture.

Dans MS-STREAM, le client utilise donc les groupes d'images redondants de basse qualité dans l'éventualité où ceux en bonne qualité ne sont pas reçus à temps. La surcharge subie par le réseau en bande passante dépend alors de la qualité des vidéos. En moyenne, nous observons moins de 10% d'augmentation de la bande passante utilisée lors de nos évaluations. De plus, la génération et agrégation des sous-segments a une empreinte minimale puisqu'il suffit d'assembler des groupes d'images déjà encodés en différentes qualités par ailleurs.

4 Muslin

Les services de streaming dépendent de larges réseaux de serveurs pour héberger le contenu vidéo. Les utilisateurs sont automatiquement redirigés vers le serveur le plus proche d'eux afin de mitiger les saturations et atteindre un meilleur débit. Cependant, si un large nombre d'utilisateurs situés dans la même région visionnent simultanément un flux vidéo, le serveur le plus proche peut rapidement être surchargé.

Muslin est une solution de streaming vidéo fournissant une *qualité d'expérience* haute et équitable aux utilisateurs, nécessitant une infrastructure moindre que les solutions actuelles. **Muslin** implémente MS-STREAM pour la livraison du contenu afin d'agréger les bandes passantes. **Muslin** utilise des retours périodiques automatisés des lecteurs vidéo des clients pendant les sessions de streaming ainsi qu'un score de classement pour provisionner et affecter dynamiquement les serveurs selon de multiples critères. Cela permet d'ajuster l'échelle de l'infrastructure en temps réel en fonction du besoin constaté et donc réduire les coûts.

Comme montré sur la Figure 4, le module de provisionnement ajuste dynamiquement le nombre de serveurs en fonction des besoins constatés et de l'estimation de la bande passante nécessaire.

Le module de sélection affecte des serveurs aux clients en fonction de plusieurs critères, tels que la distance, bande passante et charge, agrégés dans un score de classement RS_{sc} (voir Figure 5a). Comme illustré sur la Figure 5b, le serveur le plus proche n'est pas toujours le plus pertinent à affecter aux utilisateurs.

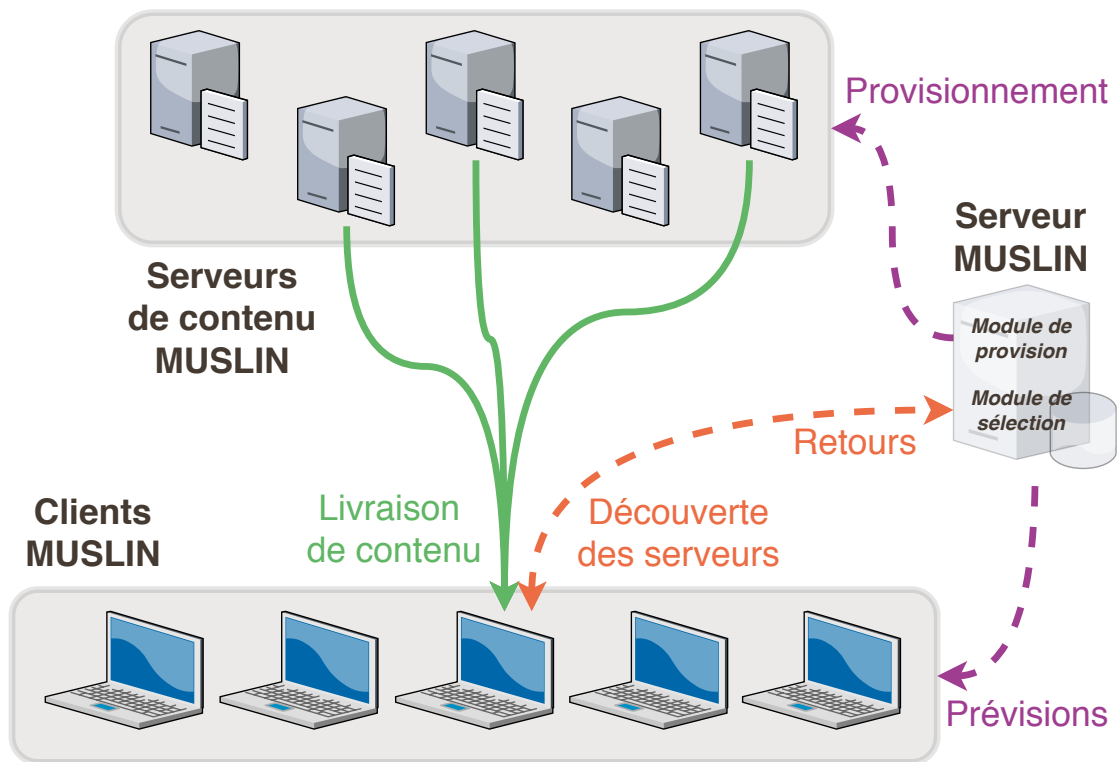
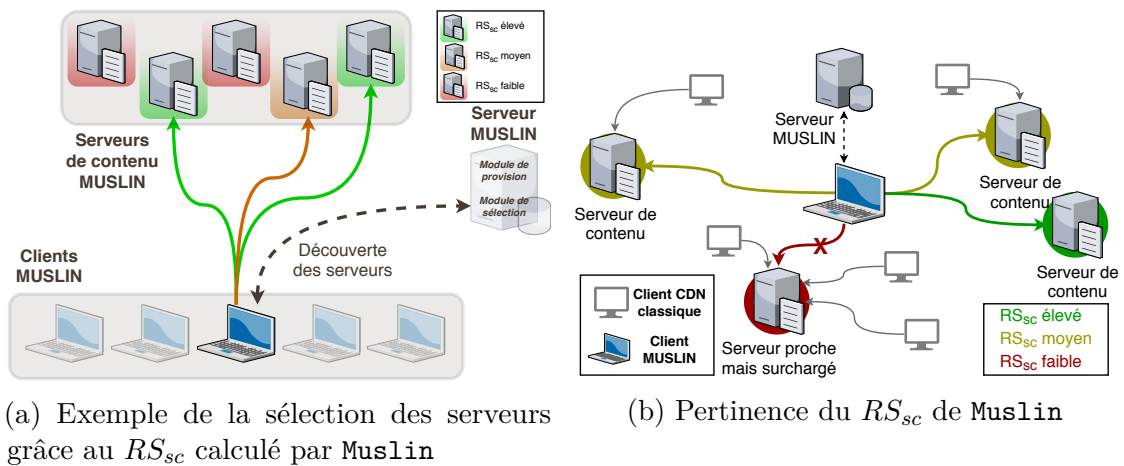


FIGURE 4 – Vue d'ensemble de Muslin

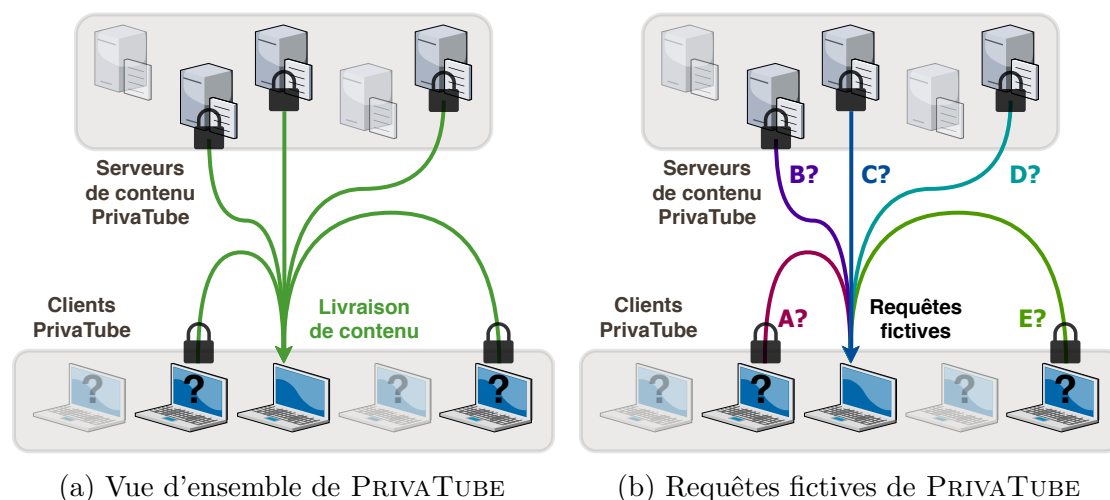


Nous avons utilisé Muslin pour rejouer un flux d'une journée couvrant un événement de jeux vidéos, avec plusieurs centaines de clients et en testant des configurations différentes. Les résultats montrent que notre approche surpasse les méthodes classiques en améliorant la *qualité d'expérience* et l'équité entre

utilisateurs (suppression totale des coupures, amélioration de la qualité vidéo et diminution des fluctuations), tout en nécessitant moins de serveurs (environ -18%).

5 PRIVATUBE

Transmettre des flux vidéo de manière fiable et à large échelle requiert une grande plateforme de diffusion avec de nombreux serveurs. Cependant, les historiques d'accès peuvent révéler des informations sensibles, et les plateformes d'hébergement sont connues pour exploiter les données personnelles. Il est donc nécessaire de protéger les intérêts des utilisateurs pour concevoir une nouvelle génération de services de streaming.



PRIVATUBE est un système de streaming vidéo pragmatique fournissant une bonne *qualité d'expérience* à ses utilisateurs tout en protégeant leur *vie privée*. PRIVATUBE étend MS-STREAM pour améliorer la *qualité d'expérience*, réduire la charge sur les serveurs et le coût de l'infrastructure en permettant aux clients de récupérer les contenus à la fois depuis les serveurs centraux et depuis les pairs ayant regardé le même contenu précédemment (voir Figure 6a). PRIVATUBE protège la *vie privée* des utilisateurs en chiffrant tous les flux dans des environnements d'exécution de confiance Intel SGX, à la fois côté client et serveur (voir Figure 7). De plus, des requêtes fictives permettent de brouiller les pistes (voir Figure 6b).

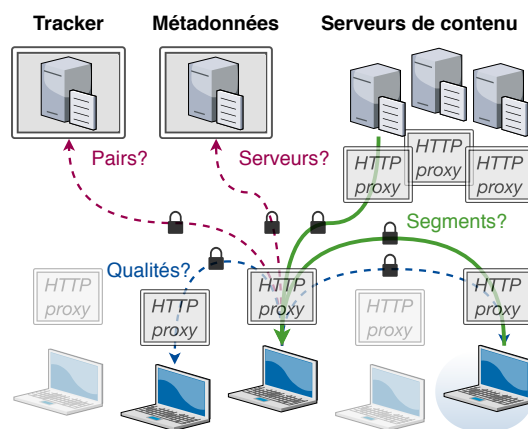


FIGURE 7 – Architecture de PRIVATUBE

En effet, PRIVATUBE permet de paramétrer la découverte probabiliste du lien entre un utilisateur et une vidéo à hauteur d'un pourcentage δ . Ces requêtes fictives ont un surcoût qui est exploité par le système en pré-provisionnement du contenu chez les pairs afin d'améliorer la disponibilité et le passage à l'échelle. Cela permet ainsi d'améliorer la *qualité d'expérience* grâce à l'agrégation des bandes passantes depuis plusieurs sources, notamment pour les vidéos moins populaires.

Nous avons implémenté PRIVATUBE et l'avons déployé sur un réseau de 14 machines pour évaluer ses performances et son comportement. Nous avons également conduit des simulations à grande échelle sur des jeux de données réels d'historiques d'accès à des vidéos. Nos résultats démontrent que PRIVATUBE offre un anonymat quasi-total aux utilisateurs tout en proposant une meilleure *qualité d'expérience* que les systèmes actuels. La durée de téléchargements des segments est 2 à 15 fois plus rapide, la qualité vidéo entre 10% et 300% plus élevée, pour un surcoût de charge serveur de 17% et un délai de démarrage supplémentaire de seulement 40ms.

6 PProx

Les plateformes de streaming vidéo (telles que YouTube, Vimeo ou Dailymotion), proposent des recommandations de contenus aux utilisateurs afin de les conserver sur leur site ou application. Pour cela, elles peuvent soit établir des profils d'intérêts

pour les utilisateurs, soit dépendre de services de recommandations externes. Le calcul de ces recommandations est toujours basé sur l'historique de navigation, et parfois sur des données entrées par les utilisateurs. Cela pose donc des menaces à la *vie privée*, puisque (i) les fournisseurs de service collectent des données personnelles, (ii) un attaquant peut intercepter les recommandations et déduire des informations privées sur l'utilisateur, et (iii) des plateformes malveillantes peuvent cibler des utilisateurs spécifiques avec de la publicité pour générer des revenus, au lieu de les segmenter par groupes d'intérêts.

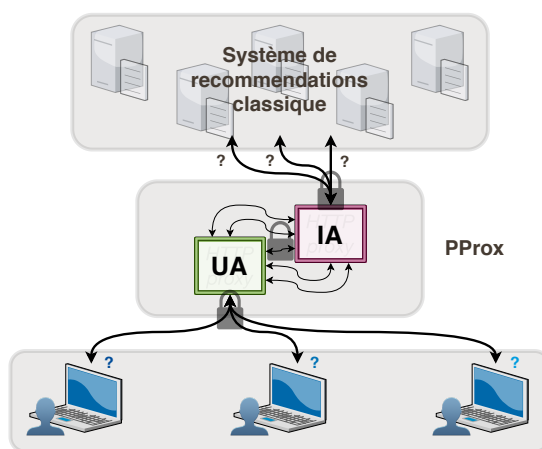
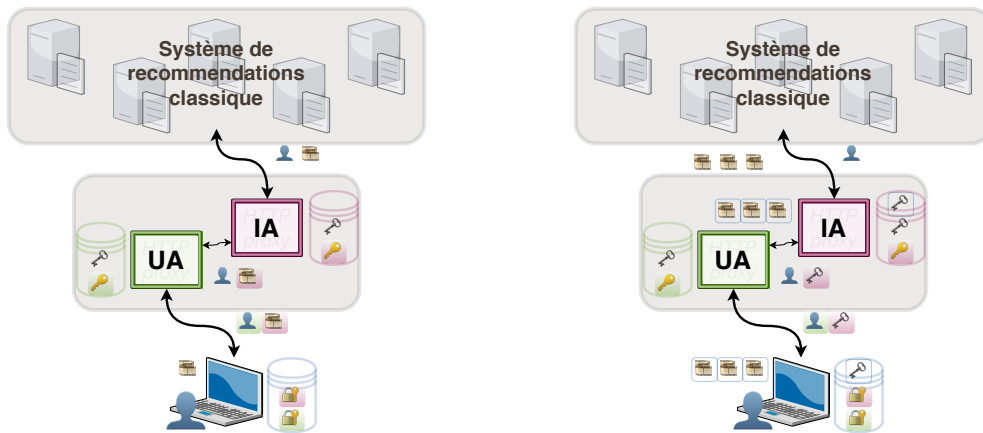


FIGURE 8 – Vue d'ensemble de PProx

PProx est une solution pragmatique permettant de fournir un service de recommandations aux utilisateurs des plateformes de streaming tout en préservant leur *vie privée*, en garantissant un anonymat total. PProx permet une bonne *qualité d'expérience* puisqu'il n'impacte pas la précision ou la nature des recommandations, et peut être déployé avec des contraintes minimales. Il dépend d'un système de double proxy dans des enclaves Intel SGX, situé entre l'utilisateur et le service de recommandations, qui chiffre et anonymise les requêtes à la volée de manière transparente. Il mélange également les requêtes des différents clients afin de casser définitivement le lien entre les utilisateurs et les contenus qu'ils visionnent ou reçoivent comme recommandation (voir Figure 8). Ce principe est robuste aux attaques de type *side-channel*, et même à la compromission d'une des enclaves. PProx passe à l'échelle de manière élastique et dynamique sur un réseau de machines disposant d'enclaves Intel SGX.



(a) PProx - Insertion d'un élément

(b) PProx - Envoi des recommandations

Nous avons connecté PProx avec le système de recommandations intégré dans Harness et l'avons évalué sur un cluster de 27 machines. Les résultats démontrent la capacité de PProx à gérer un grand nombre de requêtes avec une faible latence (moins de 100ms contre plusieurs secondes pour les systèmes similaires actuels), permettant d'atteindre la charge maximale supportée par le système de recommandations avec un surcoût acceptable (seulement 30% à 50% de nœuds en plus).

7 Conclusion

Le streaming vidéo évolue très rapidement et rencontre de nombreux défis techniques et technologiques. Nous croyons que notre travail prouve que sécurité et protection de la vie privée des utilisateurs ne rime plus avec faibles performances et basse qualité d'expérience. De nombreuses solutions pragmatiques peuvent être développées dans l'industrie en se basant sur nos contributions.

Nous espérons que généraliser Muslin, PRIVATUBE et PProx pourra permettre à une nouvelle génération de services de streaming vidéo respectant la vie privée de voir le jour.