

Complexité de Communication

Extrait d'un chapitre écrit par :
Iordanis Kerenidis, Frédéric Magniez et David Xiao
Adapté par Samuel Guerin

23 novembre 2014

Résumé

Imaginons un scénario où deux protagonistes (ou joueurs), Alice et Bob, souhaitent calculer de manière collaborative une fonction $f(x, y)$ alors que x est détenue par Alice et y par Bob. La complexité de communication de la fonction f est le nombre de bits qu'Alice et Bob doivent s'échanger afin de pouvoir calculer $f(x, y)$. Les autres ressources, telles que temps et espace mémoire, sont ici ignorées afin de se focaliser uniquement sur la communication requise.

Plus formellement, Alice et Bob s'échangent des bits à tour de rôle en suivant un protocole, qui définit à chacun instant qui doit envoyer le prochain bit, ou si le protocole est terminé. Lorsque le protocole est terminé, la valeur de $f(x, y)$ doit alors être connue des deux joueurs, ou bien d'un seul des deux joueurs dans le contexte de communication unilatérale (où seule Alice envoie un message à Bob). Ce scénario modélise une communication parfaite, où chaque message arrive à son destinataire sans être perdu ou modifié. Ce sera le cas de tout ce chapitre.

L'étude de la complexité de communication permet d'unifier des preuves de bornes inférieures dans plusieurs domaines dont les automates, les machines de Turing, les circuits VLSI, les arbres de décision, les branching programs. Nous appliquerons quand à nous ces résultats aux algorithmes de streaming

Toutes les notions et résultats présentés se généralisent aux relations, mais nous ne considérerons que les fonctions.

Nous regroupons ci-dessous les principales fonctions dont la complexité de communication sera étudiée. Notons dorénavant $[n] = \{1, 2, \dots, n\}$ pour tout entier n .

$$\begin{aligned} \text{INDEX}_n : \quad & \{0, 1\}^n \times [n] \rightarrow \{0, 1\} \\ & (x, i) \mapsto x_i. \\ \text{EQUALITY}_n : \quad & \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\} \\ & (x, y) \mapsto \begin{cases} 1, & \text{si } x = y; \\ 0, & \text{sinon.} \end{cases} \\ \text{DISJOINTNESS}_n : \quad & \{0, 1\}^n \times \{0, 1\}^n \rightarrow \{0, 1\} \\ & (x, y) \mapsto \begin{cases} 1, & \text{si } x_i = y_i = 1 \text{ pour un } i \in [n]; \\ 0, & \text{sinon.} \end{cases} \end{aligned}$$

Enfin, dans tout ce chapitre les logarithmes seront pris en base 2.

1 Généralités

1.1 Cas déterministe

Lorsque l'interaction entre les deux joueurs est permise, la sortie du protocole doit pouvoir être calculée par chacun des deux joueurs. Un protocole est alors défini comme suit.

Définition 1.1 (Protocole déterministe (cas général)). *Soient X, Y deux ensembles finis. Un protocole déterministe π entre deux joueurs, Alice (A) et Bob (B), est la description d'une conversation entre Alice et Bob comme suit. Pour toute entrée $x \in X$ d'Alice et entrée $y \in Y$ de Bob :*

1. Alice envoie un premier message $m_1(x)$ de ℓ_1 bits à Bob ;
2. Bob répond par un deuxième message $m_2(y, m_1)$ de $\ell_2(m_1)$ bits à Alice ;
3. Et ainsi de suite tant que le protocole n Pour $i \geq 1$, Alice envoie $m_{2i-1}(x, m_1, m_2, \dots, m_{2i-2})$ de $\ell_{2i-1}(m_1, m_2, \dots, m_{2i-2})$ bits, et Bob $m_{2i}(y, m_1, m_2, \dots, m_{2i-1})$ de $\ell_{2i}(m_1, m_2, \dots, m_{2i-1})$ bits.
4. La fin du protocole est décidée simultanément par les deux joueurs en fonction des messages envoyés m_1, m_2, \dots ainsi que de x pour Alice, et de y pour Bob.
5. A la fin du protocole, Alice et Bob peuvent calculer une même sortie $out(x, y)$, qui dépend des messages échangés ainsi que de x pour Alice, et de y pour Bob.

La complexité de communication $C(\pi, x)$ de π sur l'entrée x est la taille en bits du message $m(x)$ d'Alice. La complexité de communication $C(\pi)$ de π est la valeur maximale de sa complexité de communication pour toutes les entrées $x \in X$. Remarquons que la taille ℓ du message d'Alice est décidée par le protocole, et non pas par l'entrée x d'Alice. En particulier, des scénarios où Bob pourrait tirer de l'information de la "non-réception" d'un message d'Alice, ou encore de la taille du message d'Alice, ne sont pas autorisés. On définit alors la notion de complexité de la communication par :

$$D(f) = \min\{C(\pi) : \text{protocole } \pi \text{ déterministe calculant } f\}.$$

Pour toute fonction $f : X \times Y \rightarrow Z$, on a déjà que

$$D(f) \leq \log(\max(|X|, |Y|)) + \log |Z|,$$

en considérant le protocole qui envoie la valeur de x , ou de y , à l'autre joueur. Celui ci calculant alors la valeur de la fonction associée, qu'il renvoie au premier joueur.

1.2 Cas probabiliste

Lorsque les joueurs, Alice et Bob, ont accès à une source aléatoire, le protocole devient probabiliste : son exécution, sa transcription et sa sortie peuvent dépendre des choix aléatoires d'Alice et Bob. La notion de calcul d'une fonction f par un protocole est alors relâchée comme pour les algorithmes en tolérant probabilité d'erreur au plus ε , pour $0 < \varepsilon < 1/2$, c'est-à-dire en imposant que la sortie du protocole soit correcte avec probabilité au moins $(1 - \varepsilon)$, pour chaque entrée (x, y) . Cette probabilité d'erreur, pour une entrée (x, y) fixée, ne dépend donc que des choix aléatoires réalisés par les joueurs.

La source aléatoire peut être de deux natures différentes :

- Soit elle est propre à chaque joueur, c'est-à-dire les bits aléatoires r_A d'Alice sont indépendants des bits aléatoires r_B de Bob. On parle alors d'aléa privé ;
- Soit elle est commune aux deux joueurs, c'est-à-dire que les bits aléatoires d'Alice sont identiques à ceux de Bob ($r_A = r_B$). On parle alors d'aléa public.

Etant donnée une fonction f , les notions de complexités des protocoles probabilistes sont alors définies comme suit :

- Complexité de communication probabiliste à aléa privé :

$$R_\varepsilon(f) = \min \left\{ C(\pi) : \begin{array}{l} \text{protocole probabiliste } \pi \text{ calculant } f \\ \text{avec aléa privé et erreur bornée } \varepsilon \end{array} \right\};$$

- Complexité de communication probabiliste à aléa public :

$$R_\varepsilon^{\text{pub}}(f) = \min \left\{ C(\pi) : \begin{array}{l} \text{protocole probabiliste } \pi \text{ calculant } f \\ \text{avec aléa privé et public, et erreur bornée } \varepsilon \end{array} \right\}.$$

Disposant d'un aléa publique, les joueurs peuvent décider de ne garder chacun que un bit aléatoire sur deux. Ils disposent alors d'une source aléatoire privée. De même il peuvent décider d'ignorer la source. On a donc :

$$R_\varepsilon^{\text{pub}}(f) \leq R_\varepsilon(f) \leq D(f).$$

Exemple 1.2. Alors que $D(\text{EQUALITY}_n) \leq n + 1$ et on verra plus tard que $D(\text{EQUALITY}_n) = n + 1$, les protocoles 1 et 2 permettent d'établir que

$$R(\text{EQUALITY}_n) \leq \lceil \log(n/\varepsilon) \rceil + 1 \quad \text{et} \quad R_\varepsilon^{\text{pub}}(\text{EQUALITY}_n) \leq \lceil \log(1/\varepsilon) \rceil.$$

Ces protocoles renvoient toujours ' $x = y$ ' lorsque $x = y$. Leurs complexités vérifient l'inégalité ci-dessus. Lorsque $x \neq y$, le protocole 1 se trompe avec probabilité au plus ε car le polynôme $P = \sum_{i=1}^n (x_i - y_i)a^{n-i}$ s'annule alors en au plus n points (car lorsque p est premier, l'anneau \mathbb{Z}_p est un corps), et que a est pris au hasard dans un ensemble de taille au moins n/ε . Pour le protocole 2, observons que si $w \in \{0,1\}^n$ et $w \neq 0^n$, alors $\bigoplus_{i=1}^n w_i r_i = 0$ avec probabilité $1/2$, lorsque r est choisi uniformément au hasard dans $\{0,1\}^n$. Par conséquent, $w := (x_i \oplus y_i)_{i=1,2,\dots,n}$ satisfait $w \neq 0^n$ lorsque $x \neq y$. De plus $u = v$ si et seulement si $\bigoplus_{i=1}^n w_i r_i^t = 0$, pour $t = 1, 2, \dots, k$. Donc le protocole 2 se trompe avec probabilité $1/2^k \leq \varepsilon$.

Protocole 1 :

- Alice :

- Choisir un nombre premier p quelconque tel que $n/\varepsilon \leq p \leq 2n/\varepsilon$
- Choisir uniformément au hasard un entier a tel que $0 \leq a < p$
- Calculer $u = \sum_{i=1}^n x_i a^{n-i} \pmod p$
- Envoyer $m_1 := (p, a, u)$ à Bob

- Bob :

- Calculer $v = \sum_{i=1}^n y_i a^{n-i} \pmod p$
- Si $u = v$, stopper et déclarer ' $x = y$ '
- Sinon stopper et déclarer ' $x \neq y$ '

Protocole 2 :

- Alice et Bob :

— Choisir $k := \lceil 1/\varepsilon \rceil$ mots r^1, r^2, \dots, r^k de n bits uniformément au hasard

- Alice :

— Calculer $u_t = \bigoplus_{i=1}^n x_i r_i^t$, pour $t = 1, 2, \dots, k$

— Envoyer $m_1 := u$ à Bob

- Bob :

— Calculer $v_t = \bigoplus_{i=1}^n y_i r_i^t$, pour $t = 1, 2, \dots, k$

— Si $u = v$, stopper et déclarer ' $x = y$ '

— Sinon stopper et déclarer ' $x \neq y$ '

2 Communication unilatérale

Définition 2.1 (Protocole déterministe unilatéral). Soient X, Y deux ensembles finis. Un protocole déterministe unilatéral π entre deux joueurs, Alice et Bob, est la description d'une conversation unilatérale d'Alice vers Bob comme suit. Pour toute entrée $x \in X$ d'Alice et entrée $y \in Y$ de Bob :

1. Alice envoie un message $m(x)$ de ℓ bits à Bob, où $m : X \rightarrow \{0, 1\}^\ell$;
2. Bob calcule la sortie $out(m, y)$, où out est définie sur $\{0, 1\}^\ell \times Y$.

Puisque m est définie par l'entrée x , la sortie out du protocole est donc aussi une fonction de (x, y) qu'on pourra noter $out(x, y)$ au lieu de $out(m, y)$.

Comme pour les algorithmes, un protocole π calcule une fonction f si $out(x, y) = f(x, y)$ pour toutes entrées (x, y) . Etant donnée une fonction f , définissons la complexité de communication déterministe unilatérale par la quantité

$$\vec{D}(f) = \min\{C(\pi) : \text{protocole } \pi \text{ déterministe unilatéral calculant } f\}.$$

Puisqu'Alice peut toujours se contenter d'envoyer la totalité de son entrée x à Bob, on peut déjà établir que

$$\vec{D}(f) \leq \lceil \log |X| \rceil.$$

Comme nous le voyons dans l'exemple précédent, $\vec{D}(f)$ peut être beaucoup plus petite que le majorant $\lceil \log |X| \rceil$. Afin de caractériser $\vec{D}(f)$, introduisons la notion de matrice de communication M_f de f . Cette matrice est indexée par $x \in X$, pour les lignes, et par $y \in Y$, pour les colonnes, et satisfait

$$(M_f)_{x,y} = f(x, y).$$

Théorème 2.2. Soient X, Y deux ensembles finis, et f une fonction définie sur $X \times Y$.

$$\vec{D}(f) = \lceil \log(\text{nombre de lignes distinctes dans } M_f) \rceil.$$

Démonstration. Posons $k =$ nombre de lignes distinctes dans M_f , et considérons un protocole π quelconque pour f . Nous commençons par montrer qu'Alice doit pouvoir envoyer au moins k messages différents lorsque x varie, c'est-à-dire que $C(\pi) \geq \lceil \log(k) \rceil$. Pour cela, montrons que si deux lignes x et x' de M_f sont distinctes alors nécessairement $m(x) \neq m(x')$. Par contraposition, si $m(x) = m(x')$ alors pour toute entrée $y \in Y$ le protocole calcule $out(m(x), y) = out(m(x'), y)$.

Protocole 1 – Protocole unilatéral optimal

- Préparation :
 Soit L_1, L_2, \dots, L_k les lignes sans multiplicité de M_f
 ordonnées arbitrairement (ordre connu d’Alice et Bob)

- Alice :
 Soit $i \in [k]$ l’entier tel que $L_i = (M_f(x, y))_{y \in Y}$
 Envoyer $m := i$ à Bob

- Bob : Calculer $out := L_m[y]$

Puisque π calcule f , il s’en suit que $f(x, y) = f(x', y)$, pour tout $y \in Y$, et donc que les lignes x et x' de M_f sont identiques.

Pour montrer $\vec{D}(f) \leq \lceil \log k \rceil$, nous construisons un protocole π pour f tel que $C(\pi) = \lceil \log k \rceil$. Soit L l’ensemble des lignes de M_f . Par hypothèse, l’ensemble L est de taille k . Le protocole 1 calcule alors correctement f et est de complexité annoncée. \square

Une conséquence de ce résultat est la caractérisation des complexités suivantes, dont nous laissons les preuves en exercice :

Lemme 2.3.

$$\vec{D}(\text{EQUALITY}_n) = \vec{D}(\text{INDEX}_n) = \vec{D}(\text{DISJOINTNESS}_n) = n.$$

3 Communication bilatérale déterministe

La transcription $m(x, y)$ de π sur l’entrée (x, y) est la concaténation des messages échangés $m(x, y) := m_1 m_2 \dots$. La complexité (de communication) $C(\pi, x, y)$ de π sur l’entrée (x, y) est la taille en bits de la transcription correspondante. La complexité (de communication) $C(\pi)$ de π est la valeur maximale de sa complexité de communication pour toutes les entrées $(x, y) \in X \times Y$.

Une propriété fondamentale des protocoles de communication est que les entrées (x, y) menant à la même transcription forment un rectangle : Un *rectangle* est un ensemble R tel que pour tout $(x, y) \in R$ et $(x', y') \in R$, alors $(x, y') \in R$ et $(x', y) \in R$. Un corollaire est que la sortie d’un protocole ne dépend en fait que de sa transcription. Il est donc possible à quiconque qui écoute la communication entre Alice et Bob de calculer la fonction sans même connaître x ou y !

Lemme 3.1. *Soit π un protocole déterministe sur $X \times Y$. Soit une paire d’entrées (x, y) et (x', y') de π ayant même transcription t . Alors (x', y) et (x, y') ont aussi même transcription, de plus cette transcription est t .*

En particulier, les sorties de π sur les entrées (x, y) , (x', y) , (x, y') et (x', y') sont identiques.

Preuve (esquisse). Considérons deux paires d’entrées (x, y) et (x', y') ayant même transcription $t = m(x, y) = m(x', y')$. Posons $t = t_1 t_2 \dots$ avec $t_i = m_i(x, y)$. Par récurrence sur i , on montre que le i -ème message de π sur (x', y) est t_i , et de façon similaire sur (x, y') .

Puisque la sortie $out(x, y)$ ne dépend que de la transcription et x , ou de la transcription et y , on en déduit que $out(x, y) = out(x', y)$, puis que $out(x', y) = out(x', y')$, soit au final que $out(x, y) = out(x', y')$. \square

Cette propriété simple permet de montrer des bornes inférieures importantes, comme le théorème suivant.

Théorème 3.2. *Soient X, Y deux ensembles finis. Soient f une fonction définie sur $X \times Y$ et $S \subseteq X \times Y$ tels que*

1. *f est constante sur S ;*
2. *tout $(x, y), (x', y') \in S$ satisfait $f(x, y) \neq f(x', y)$ ou $f(x, y) \neq f(x, y')$.*

Alors

$$D(f) \geq \lceil \log(1 + |S|) \rceil.$$

Démonstration. Soit $z = f(x, y)$, pour un élément $(x, y) \in S$ quelconque, la valeur que f prend sur S . La contraposée du Lemme 3.1 et l'hypothèse (2) impliquent que chaque paire d'éléments $(x, y), (x', y') \in S$ a nécessairement une transcription différente. Donc il faut au moins $|S|$ transcriptions différentes pour avoir un protocole qui calcule correctement f .

Puisque la fonction f n'est pas constante, le Lemme 3.1 implique aussi qu'il faut au moins une transcription de plus pour calculer f , d'où le résultat. \square

Le théorème 3.2 permet de montrer que

Lemme 3.3.

$$D(\text{EQUALITY}_n) = D(\text{DISJOINTNESS}_n) = n + 1.$$

Démonstration. La majoration est obtenue par le protocole trivial du Lemme ???. Pour obtenir la minoration il suffit de considérer pour EQUALITY_n l'ensemble $S = \{(x, x) : x \in \{0, 1\}^n\}$, et pour DISJOINTNESS_n l'ensemble $S = \{(x, \bar{x}) : x \in \{0, 1\}^n\}$, où \bar{x} est obtenu en prenant la négation de tous les bits de x . Dans les deux cas, l'ensemble S est de taille 2^n et satisfait les hypothèses du théorème 3.2. Par conséquent, la quantité $\lceil \log(1 + 2^n) \rceil = n + 1$ est un minorant de chacune des complexités. \square

Cet autre lemme permet d'établir une borne inférieure sur $D(f)$ à partir de propriétés algébriques de M_f . Notons $\delta_z(f)$ la fonction qui prend la valeur 1 en (x, y) si $z = f(x, y)$ et la valeur 0 sinon.

Théorème 3.4. *Soient X, Y deux ensembles finis, et $f : X \times Y \mapsto Z$.*

$$D(f) \geq \left\lceil \log \left(\sum_{z \in Z} \text{rang}(M_{\delta_z(f)}) \right) \right\rceil.$$

Démonstration. Fixons un protocole π pour f . Etant donnée une transcription t possible de π , appelons $R_t = \{(x, y) \in X \times Y : m(x, y) = t\}$. Nous interpréterons aussi R_t comme la matrice booléenne indexée par $x \in X$ et $y \in Y$, et prenant la valeur 1 en $(x, y) \in R_t$ et 0 sinon. D'après le lemme 3.1, nous savons que R_t est un rectangle, et que toute entrée de R_t mène vers une même sortie $out(t)$. Pour $z \in Z$, définissons $S_z := \{t : out(t) = z\}$. Alors

$$M_{\delta_z(f)} = \sum_{t \in S_z} R_t.$$

Le théorème du rang nous dit alors que

$$\text{rang}(M_{\delta_z(f)}) \leq \sum_{t \in S_z} \text{rang}(R_t) = |S_z|,$$

où l'égalité provient du fait que chaque R_t est non nul, et donc de rang au moins 1, pour $t \in S_z$.

La preuve se conclut en observant que $\sum_{z \in Z} |S_z|$ représente le nombre de transcriptions différentes de π peut transcrire lorsque son entrée (x, y) parcourt $X \times Y$. Par conséquent sa complexité est au moins

$$C(\pi) \geq \left\lceil \log \left(\sum_{z \in Z} |S_z| \right) \right\rceil.$$

□

Pour le problème EQUALITY_n , le théorème 3.4 nous permet de redémontrer que

$$D(\text{EQUALITY}_n) \geq n + 1.$$

En effet, observons que $M_{\delta_1(f)}$ est la matrice identité de taille $2^n \times 2^n$, et est donc de rang n . De plus $M_{\delta_0(f)}$ n'est pas nulle, et donc $D(\text{EQUALITY}_n) > n$, ou encore $D(\text{EQUALITY}_n) \geq n + 1$ puisque $D(\text{EQUALITY}_n)$ est une valeur entière.

Exemple 3.5. La valeur médiane d'une famille z_1, z_2, \dots d'éléments de $[n]$ est définie par

$$\text{médiane}(z) = \min\{v \in [n] : |\{i : z_i \geq v\}| \leq |\{i : z_i \leq v\}|\}.$$

Pour $X = Y = [n]^n$, le protocole ci dessous permet de calculer par une recherche dichotomique la valeur médiane de $(x, y) \in X \times Y$ avec une complexité de communication de $\lceil \log n \rceil (\lceil \log n \rceil + 1)$.

- Partir de $a = 1, b = n$
- Pour i allant de 1 à $\lceil \log n \rceil$
 Alice et Bob :
 - Calculer $v = \lceil (a + b)/2 \rceil$
 - Alice :
 - Calculer et envoyer $m_{2^i-1} := |\{i : x_i \geq v\}|$
 - Bob :
 - Calculer $k := m_{2^i-1} + |\{i : x_i \geq v\}|$
 - Si $k \leq n$ alors envoyer $m_{2^i} := 0$
 - Sinon envoyer $m_{2^i} := 1$
- Alice et Bob :
 - Si $m_{2^i} = 0$ alors affecter $b := v$
 - Sinon affecter $a := v$
- Fin du protocole La sortie du protocole est out $:= v$

4 Protocoles probabilistes

On procède de façon analogue pour les communications unilatérales probabilistes en définissant $\vec{R}_\varepsilon^{\text{pub}}(f)$ et $\vec{R}_\varepsilon(f)$.

Obtenir des minorants de la complexité de communication probabiliste est beaucoup plus complexe. Nous verrons plus loin une méthode générique pour cela. En attendant voici quelques premiers résultats obtenus relativement simplement. A l'aide d'arguments plus sophistiqués, il est possible de supprimer le facteur logarithmique du Lemme 4.1.

Lemme 4.1.

$$\vec{R}_{\frac{1}{3n}}(\text{INDEX}_n) \geq n, \quad \vec{R}_{\frac{1}{3}}(\text{INDEX}_n) \geq n/\log n.$$

Démonstration. La deuxième minoration se déduit de la première comme suit. Etant donné un protocole π pour INDEX_n avec erreur $1/3$, il suffit d'exécuter $\Theta(\log n)$ fois π , mais avec des choix aléatoires différents, pour obtenir une erreur au plus $1/(3n)$. En effectuant ces répétitions en parallèle le protocole reste unilatéral.

Pour montrer la première minoration, nous allons montrer qu'il est possible de dérandomiser tout protocole π pour INDEX_n d'erreur $1/(3n)$, tout en conservant sa complexité. Le lemme 2.3 permet alors de conclure. Soient r_A et r_B les bits aléatoires respectifs d'Alice et Bob. Le protocole déterministe est décrit par le protocole suivant :

- Alice choisit r_A tel que $\Pr_{r_A, r_B}[\exists j \in [n] : P(m(x, r_A), j, r_B) \neq x_j]$ est minimal.
- Alice envoie $t := m(x, r_A)$ à Bob.
- Bob calcule la valeur majoritaire de $out(t, i, r_B)$, lorsque r_B parcourt tous les choix aléatoires possibles.

Par hypothèse, pour tout $x \in \{0, 1\}^n$ et $j \in [n]$,

$$\Pr_{r_A, r_B} [out(m(x, r_A), j, r_B) \neq x_j] \leq \frac{1}{3n}.$$

Par sommation (union-bound), on obtient pour tout $x \in \{0, 1\}^n$,

$$\Pr_{r_A, r_B} [\exists j \in [n] : out(m(x, r_A), j, r_B) \neq x_j] \leq \frac{1}{3}.$$

Donc l'aléa r_A choisi par Alice et le message $t = m(x, r_A)$ envoyé à Bob satisfont nécessairement

$$\Pr_{r_B} [\exists j \in [n] : out(t, j, r_B) \neq x_j] \leq \frac{1}{3}.$$

Le message t permet de conclure pour toute entrée $i \in [n]$ de Bob. En effet, pour la valeur particulier $j = i$, l'inégalité précédente entraîne

$$\Pr_{r_B} [out(t, i, r_B) \neq x_i] \leq \frac{1}{3}.$$

Ainsi x_i est bien la valeur majoritaire de $out(t, i, r_B)$ lorsque r_B parcourt tous les choix aléatoires possibles de Bob. A noter que cette étape, ainsi que celle du choix de r_A , sont potentiellement longues mais déterministes et ne nécessitent aucune communication. \square

5 Applications

Nous allons réduire différents problèmes déjà vus précédemment à des problèmes de complexité de communication. Nous pourrons en tirer des informations très utiles, notamment des résultats d'optimalité. Rappelons ici les résultats mentionnés précédemment :

.	$\vec{D}(f)$	$D(f)$	$\vec{R}_{1/3}(f)$
INDEX_n	n	$\log(n) + 1$	$\Theta(n)$
EQUALITY_n	n	$n + 1$	$O(\log n)$
DISJOINTNESS_n	n	$n + 1$	

Considérons le problème de déterminer dans un flux de n entiers a_1, a_2, \dots de $[m]$, la plus grande multiplicité de ces entiers, notée aussi F^∞ . Nous allons plus tard relier le calcul de F^∞ en streaming à celui du calcul d'une autre fonction (INDEX) en complexité de communication.

D'abord la complexité en mémoire d'un algorithme de streaming est simplement reliée à la complexité de communication comme suit.

Lemme 5.1. *Soit f une fonction définie sur $X \times Y$ avec $X = Y = \{0, 1\}^n$. Alors tout algorithme de streaming à une passe doit utiliser une mémoire $s(n) \geq \vec{D}(f)$ s'il est déterministe, et $s(n) \geq \vec{R}_\varepsilon(f)$ s'il est probabiliste et d'erreur ε .*

De plus, tout algorithme de streaming à $p(n)$ passes doit utiliser une mémoire $s(n) \geq D(f)/(2p(n))$ s'il est déterministe, et $s(n) \geq R_\varepsilon(f)/(2p(n))$ s'il est probabiliste et d'erreur ε .

Ce lemme permet de montrer plusieurs compris mémoire et nombre de passes en streaming.

Théorème 5.2. *Tout algorithme de streaming déterministe à $p(n)$ passe, de mémoire $s(n)$, qui permet de décider si deux chaînes consécutives de n bits sont égales nécessite une mémoire $s(n) \geq n/(2p(n))$.*

Cependant, comme nous allons le voir, il est parfois nécessaire d'effectuer une réduction plus complexe, afin de transformer la fonction à calculer par un algorithme de streaming en une autre fonction plus adaptée à la complexité de communication.

Lemme 5.3. *Soit A un algorithme de streaming déterministe à une passe et de mémoire $s(n)$ permettant de calculer une valeur v telle que $|v - F^\infty| < F^\infty/3$ sur un flux de $n + 1$ de $[2n]$. Alors $\vec{D}(\text{INDEX}_n) \leq s(n)$.*

De plus, si A est probabiliste avec erreur ε , alors $\vec{R}_\varepsilon(\text{INDEX}_n) \leq s(n)$.

Démonstration. En utilisant l'algorithme de streaming A , nous allons construire un protocole π unilatéral pour INDEX_n , où Alice reçoit l'entrée $x \in \{0, 1\}^n$, et Bob l'entrée $i \in [n]$. La réduction est déterministe :

- Alice simule l'algorithme A sur le début du stream formé des entiers j tels que $x_j = 1$.
- Alice continue de simuler l'algorithme A sur des entiers $n < i' \leq 2n$ distincts jusqu'à avoir traité exactement n valeurs.
- Alice envoie la mémoire de l'algorithme A à Bob.
- Bob termine la simulation de l'algorithme A sur la fin du stream formée de l'unique entier i , et calcule la valeur v renvoyée par A .
- Si $v \geq 1/3$, alors Bob renvoie $out := 1$; sinon Bob renvoie $out := 0$

Pour prouver que le protocole a la même erreur que l'algorithme A (c'est-à-dire 0 si A est déterministe, et ε sinon), il suffit de remarquer que $F^\infty = 1$ lorsque $x_i = 0$, et que $F^\infty = 2$ sinon. Lorsque v satisfait la promesse $|v - F^\infty| < F^\infty/3$, on conclue en observant que $v < 1/3$ lorsque $x_i = 0$, et $v > 1/3$ sinon. \square

Puisque $\vec{R}_{1/3}(\text{INDEX}_n) = \Theta(n)$, nous en déduisons le théorème suivant.

Théorème 5.4. *Tout algorithme de streaming probabiliste à une passe, de mémoire $s(n)$ et erreur $\varepsilon \leq 1/3$ permettant de calculer une valeur v telle que $|v - F^\infty| < F^\infty/3$ sur un flux de $n + 1$ entiers de $[2n]$ nécessite une mémoire $s(n) = \Theta(n)$.*

Afin de généraliser ce résultat aux passes multiples, une réduction à un autre problème de communication doit être considérée. En effet, $D(\text{INDEX}_n) = O(\log n)$. Nous considérons donc maintenant une nouvelle réduction à la fonction DISJOINTNESS_n .

Lemme 5.5. *Soit A un algorithme de streaming déterministe à $p(n)$ passes et de mémoire $s(n)$ permettant de calculer une valeur v telle que $|v - F^\infty| < F^\infty/3$ sur un flux de $2n$ entiers de $[3n]$. Alors $D(\text{DISJOINTNESS}_n) \leq 2p(n)s(n)$.*

De plus, si A est probabiliste avec erreur $\varepsilon \leq 1/3$, alors $R_\varepsilon(\text{DISJOINTNESS}_n) \leq 2p(n)s(n)$.

Démonstration. La preuve est très similaire à celle du lemme 5.3. La réduction est toujours déterministe :

- Alice simule l'algorithme A sur le début du stream formé des entiers i tels que $x_i = 1$, et complété avec des entiers $n < i' \leq 2n$ distincts afin d'obtenir n valeurs.
- Bob simule l'algorithme A sur la fin du stream formé des entiers i tels que $y_i = 1$, et complété avec des entiers $2n < i' \leq 3n$ distincts afin d'obtenir n valeurs.
- Au total le stream comporte $2n$ entiers de $[3n]$, les joueurs simulant à tour de rôle l'algorithme durant ses $p(3n)$ passes.
- A la fin de la simulation, Alice et Bob calculent la valeur v renvoyée par A .
- Si $v \geq 1/3$, alors ils décident $\text{out} := 1$; sinon $\text{out} := 0$

Au total $2p(n)s(n)$ bits sont échangés. Pour prouver que le protocole a la même erreur que l'algorithme A , il suffit de remarquer que $F^\infty \leq 1$ lorsque $\text{DISJOINTNESS}_n(x, y) = 0$, et que $F^\infty = 2$ sinon. Lorsque v satisfait la promesse $|v - F^\infty| < F^\infty/3$, on conclue en observant que $v < 1/3$ lorsque $\text{DISJOINTNESS}_n(x, y) = 0$, et $v > 1/3$ sinon. \square

Nous obtenons ainsi le résultat suivant qu'il est possible de généraliser aux protocoles probabilistes.

Théorème 5.6. *Tout algorithme de streaming déterministe à $p(n)$ passes, de mémoire $s(n)$ permettant de calculer une valeur v telle que $|v - F^\infty| < F^\infty/3$ sur un flux de $2n$ entiers de $[3n]$ nécessite une mémoire $s(n) \geq (n + 1)/(2p(n))$.*