

Comment compter les moutons sans perdre la mémoire

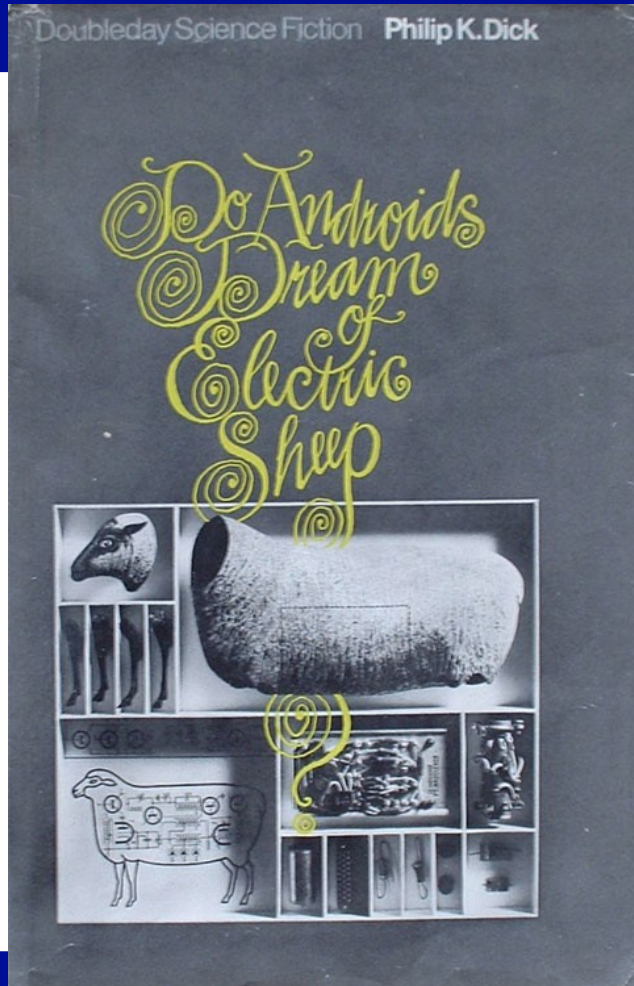


Image de catalyststuff sur Freepik

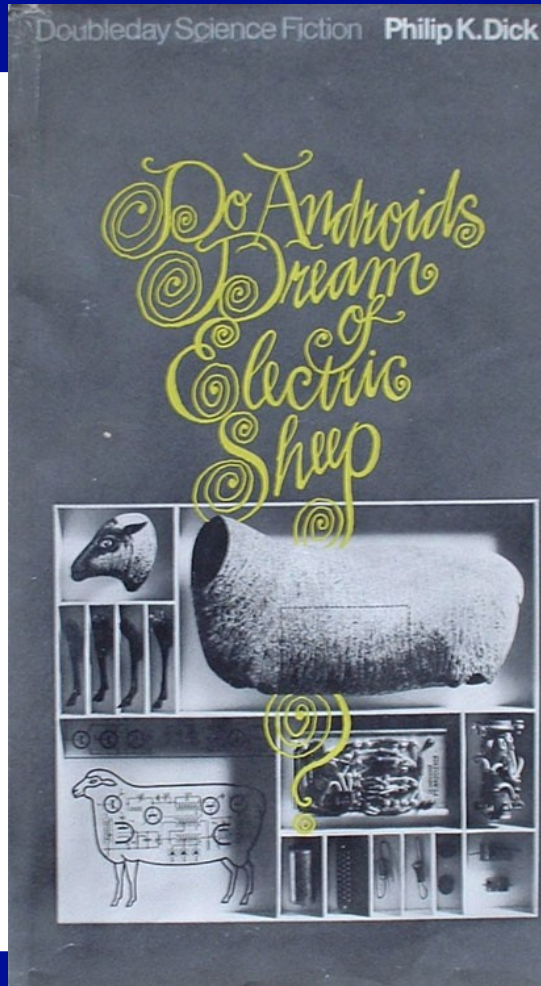
Benoît Masson



Do Androids Dream of Electric Sheep?



Do Androids Dream of Electric Sheep?



Comment compter les moutons sans perdre la mémoire

Benoît Masson



BREIZHCAMP
2049

Générique

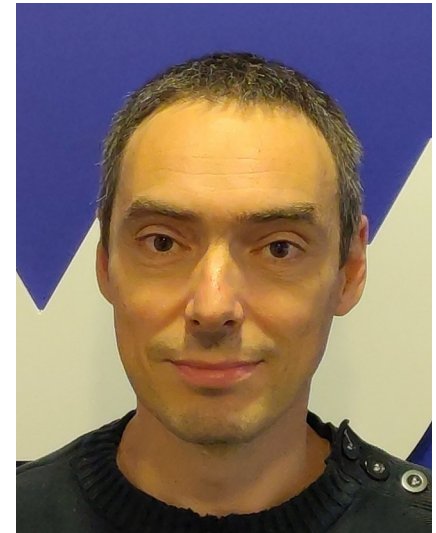
- Développeur Go depuis 2015



- Software Craftsman



@OVHcloud Rennes
Noms de Domaines depuis 2020



benoitmasson

Comptons les moutons



OVHcloud

Décompte simple (voire naïf)



Décompte simple (voire naïf)



- + performant
- mémoire non bornée
- pas de résultats intermédiaires

Décompte itératif



Décompte itératif



- + performant
- + mémoire bornée
- + visualisation en temps réel (*stream*)

Comptons les moutons (distincts)



OVHcloud

Décompte exact

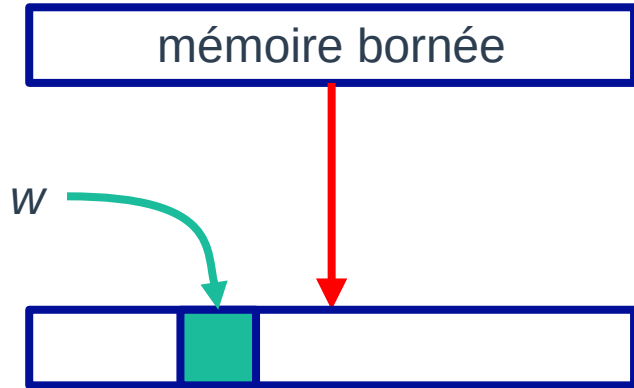


Décompte exact

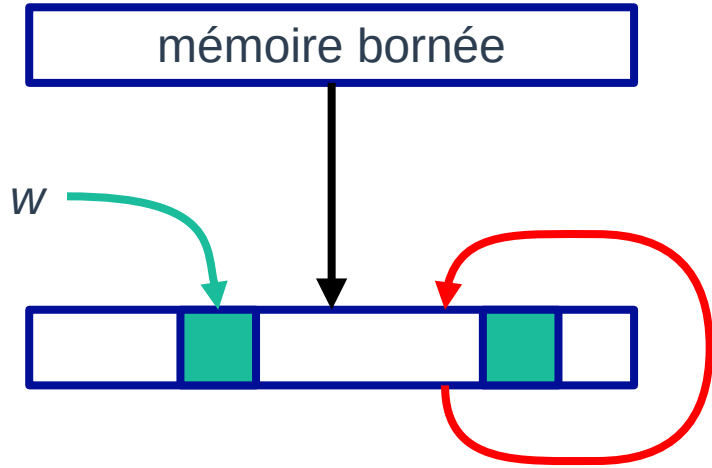


- + performant
- mémoire non bornée

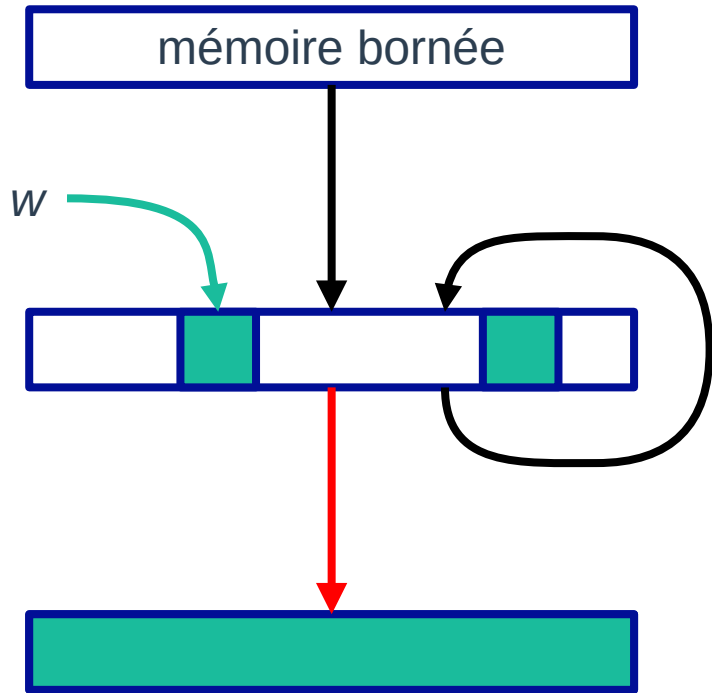
Décompte approximatif - Algorithme



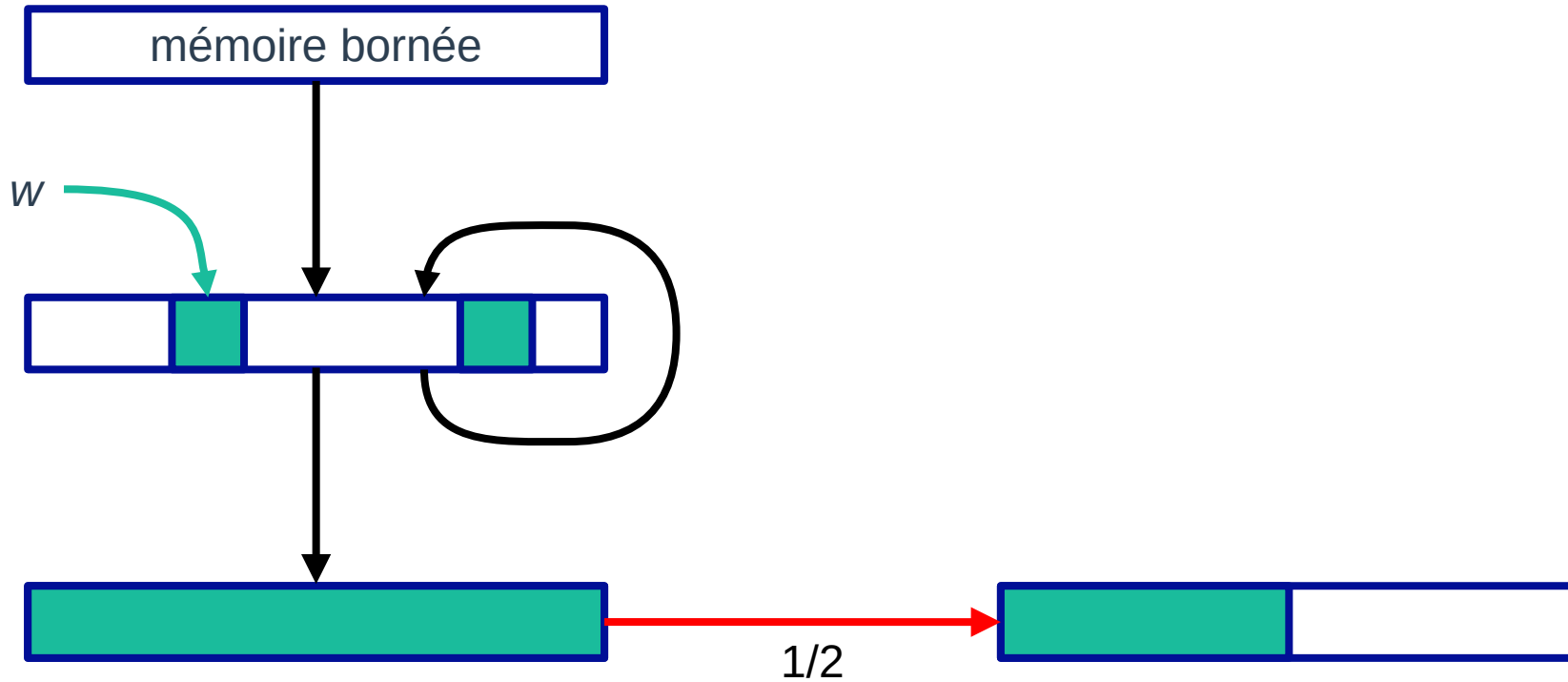
Décompte approximatif - Algorithme



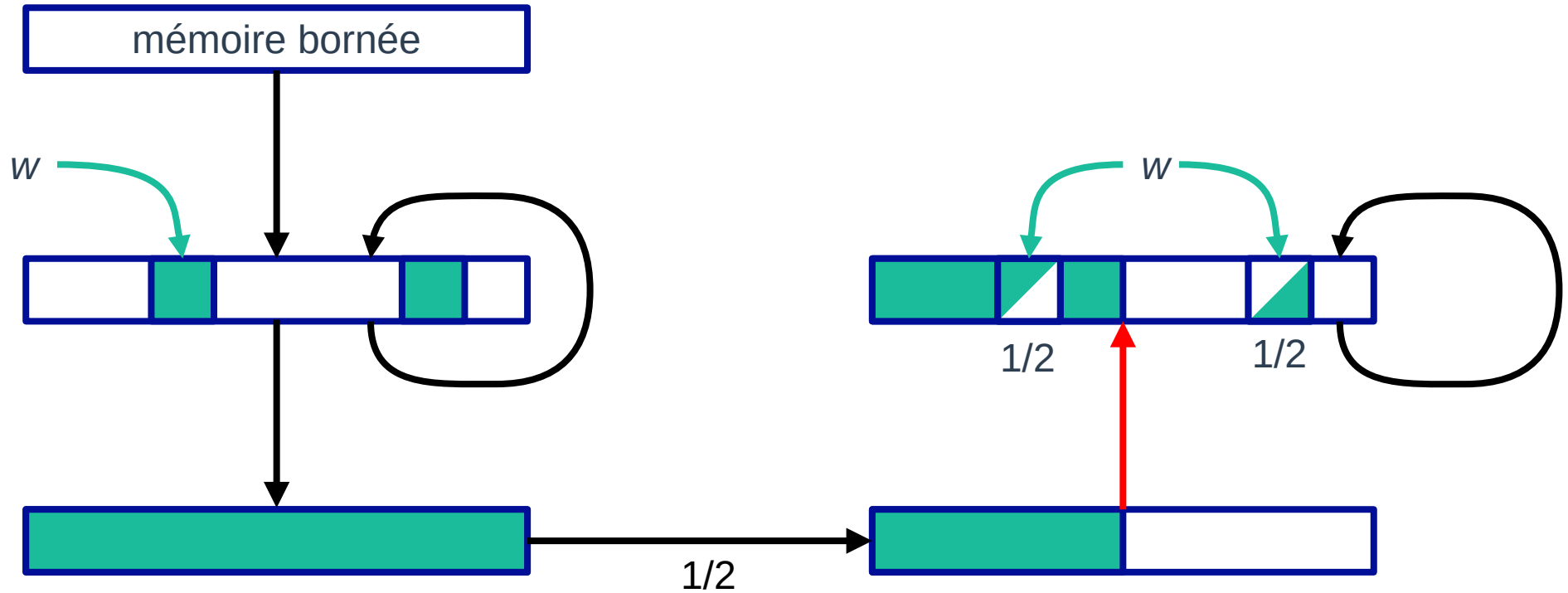
Décompte approximatif - Algorithme



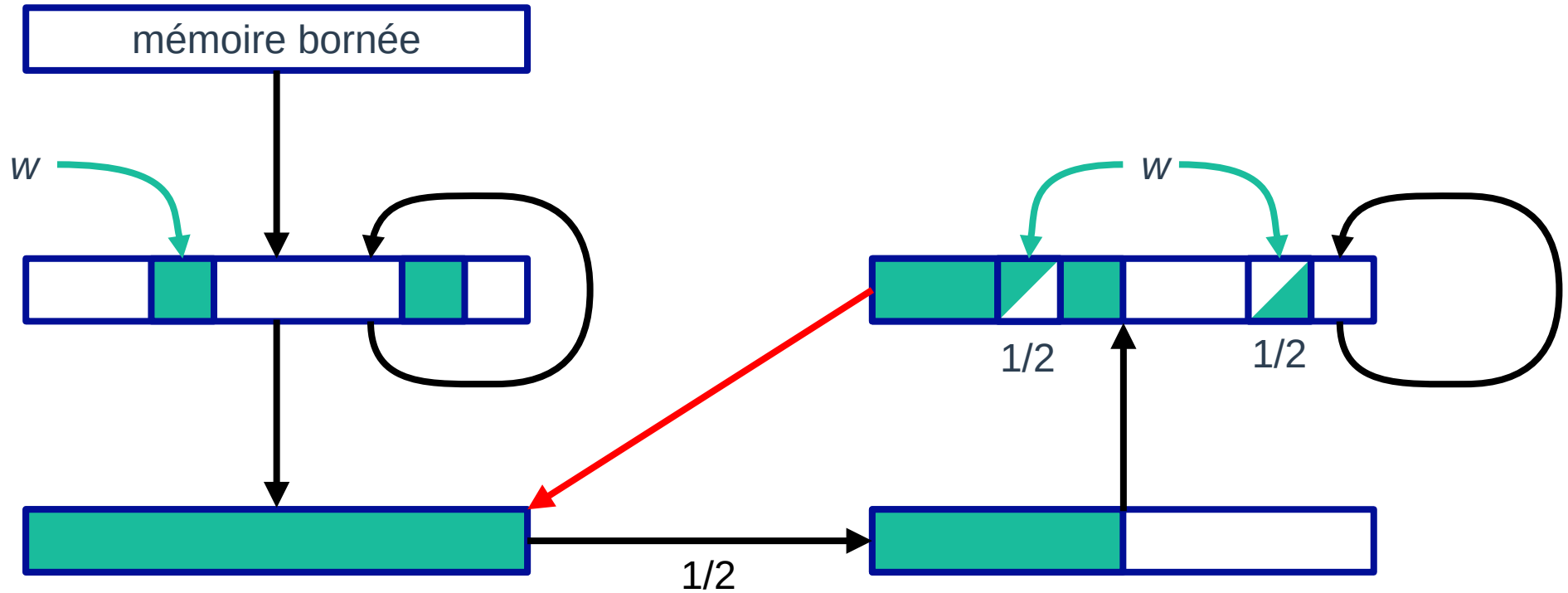
Décompte approximatif - Algorithme



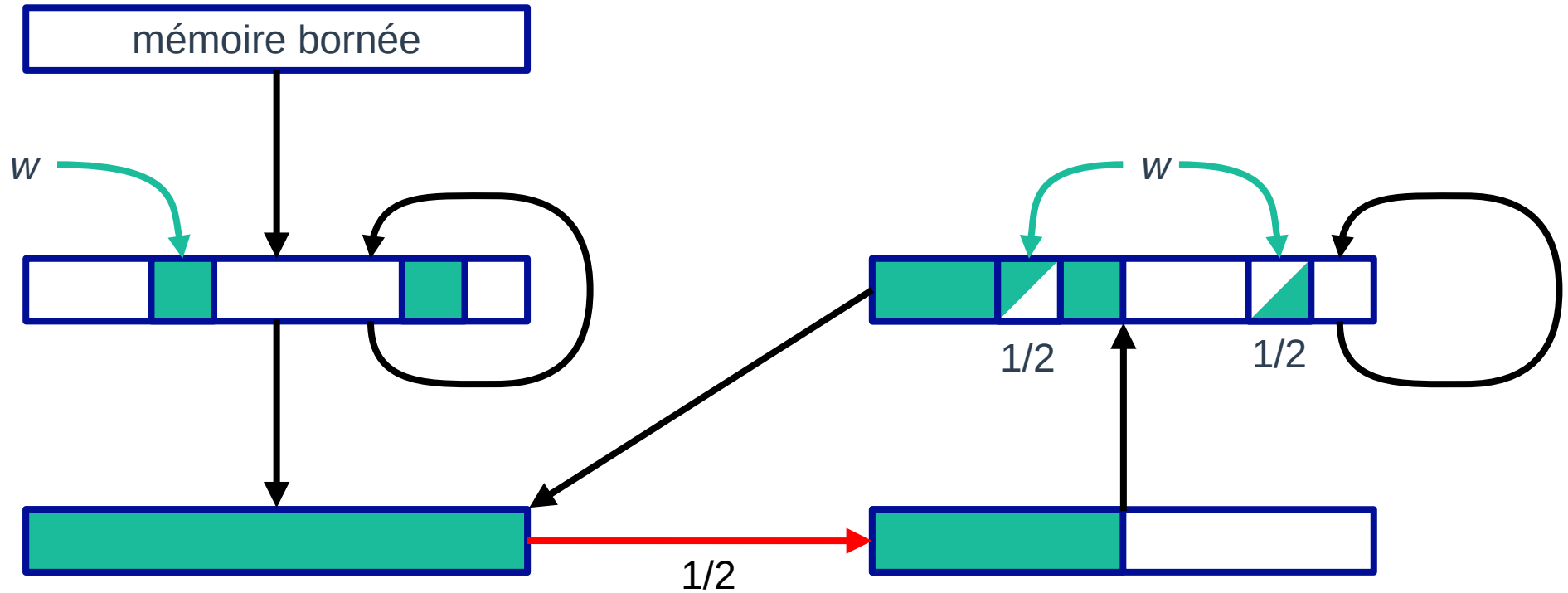
Décompte approximatif - Algorithme



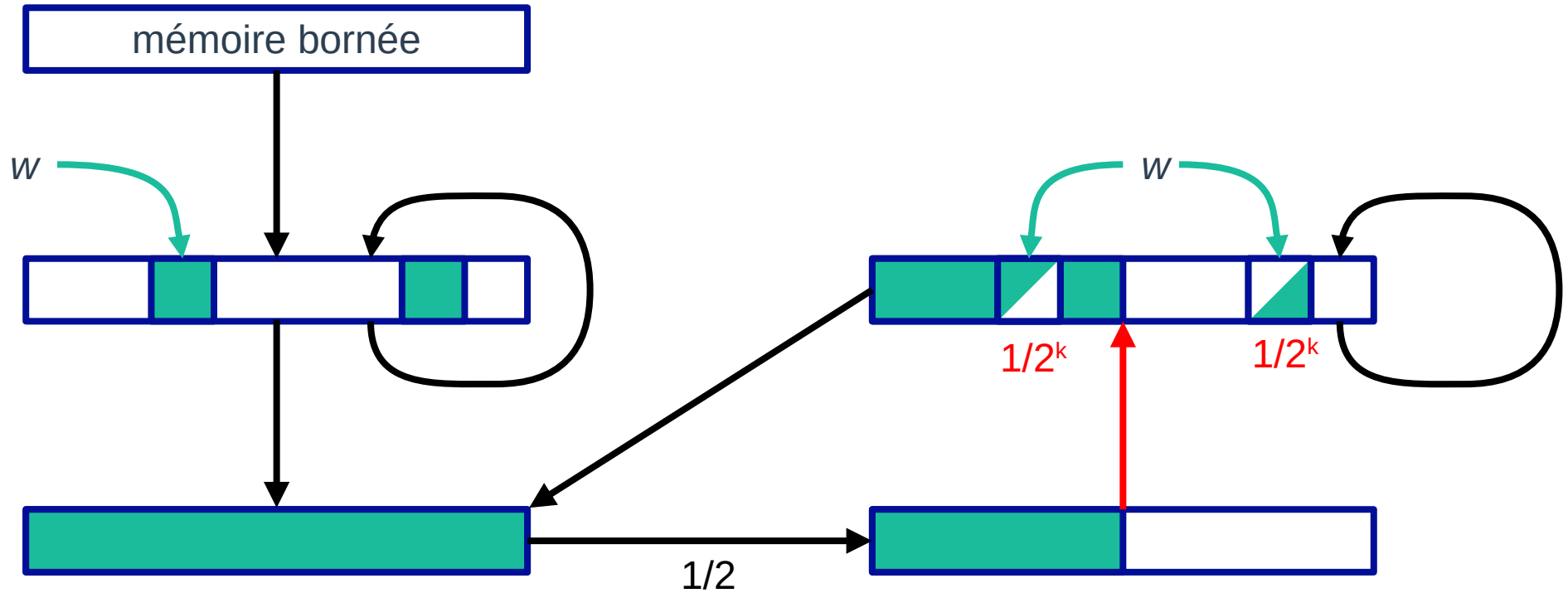
Décompte approximatif - Algorithme



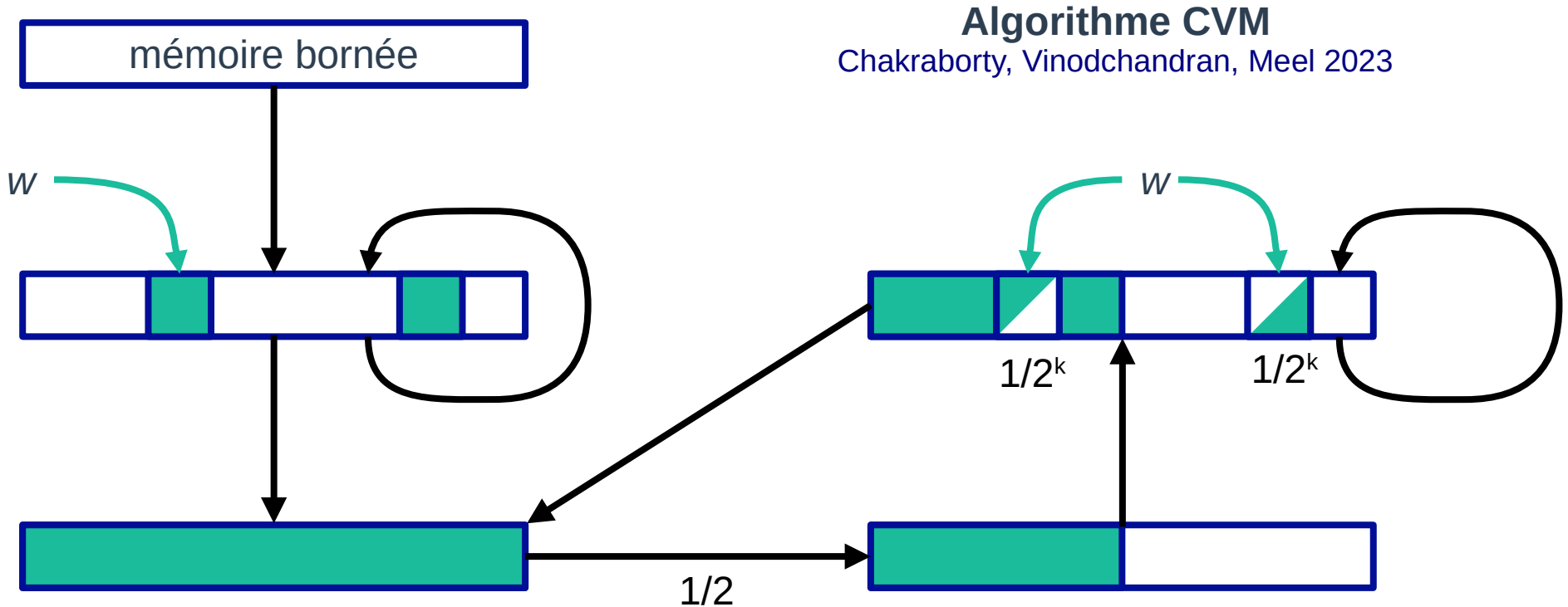
Décompte approximatif - Algorithme



Décompte approximatif - Algorithme

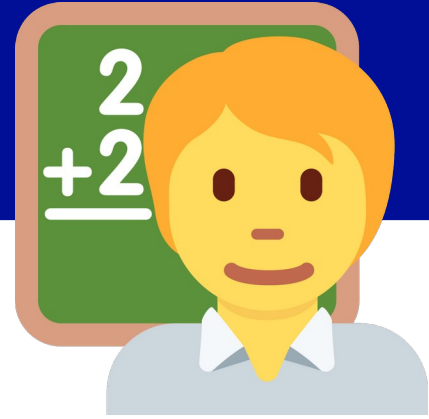


Décompte approximatif - Algorithme



Un peu de probabilités

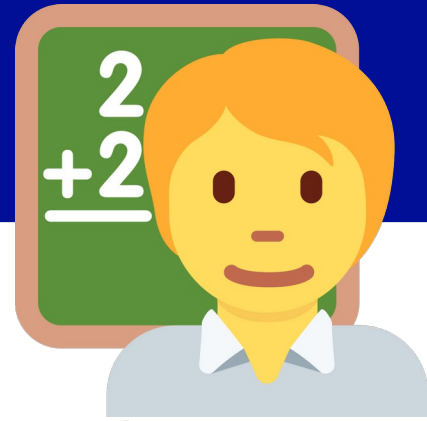
Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :



Un peu de probabilités

Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

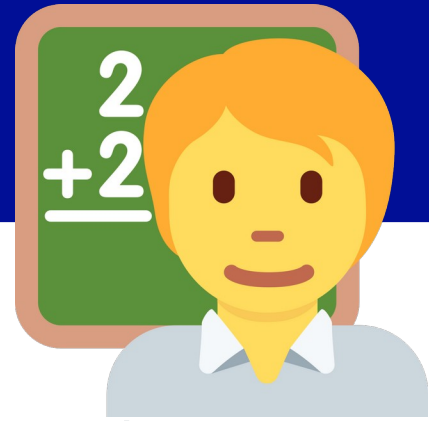
- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)



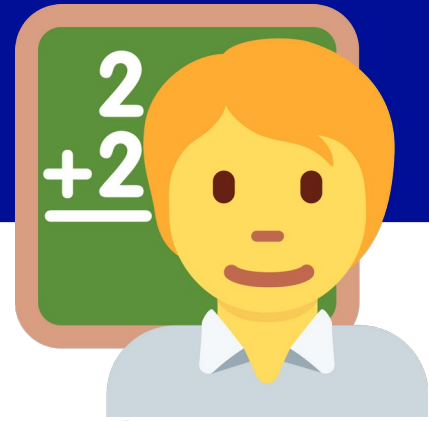
Un peu de probabilités

Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)
- si w n'apparaît qu'au second round, $P(w) = 1/2$ (ajout)



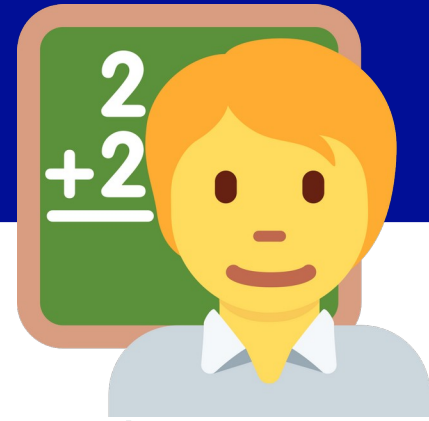
Un peu de probabilités



Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)
- si w n'apparaît qu'au second round, $P(w) = 1/2$ (ajout)
- si w apparaît dans les 2 rounds, $P(w) = 1/2$ (ajout/suppression)

Un peu de probabilités

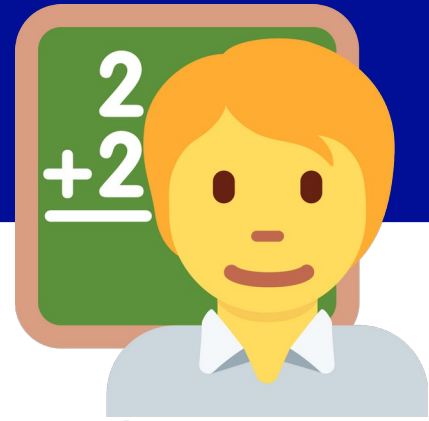


Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)
- si w n'apparaît qu'au second round, $P(w) = 1/2$ (ajout)
- si w apparaît dans les 2 rounds, $P(w) = 1/2$ (ajout/suppression)

Probabilité générique après k rounds (récurrence) : **$P(w) = 1/2^k$**

Un peu de probabilités



Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

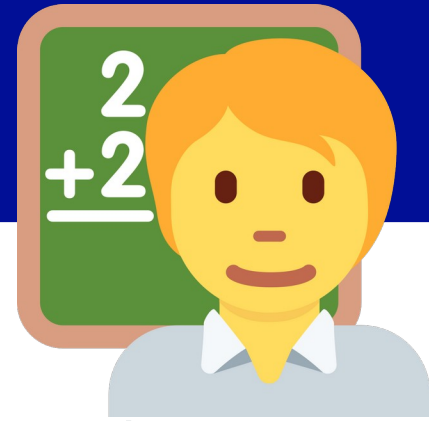
- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)
- si w n'apparaît qu'au second round, $P(w) = 1/2$ (ajout)
- si w apparaît dans les 2 rounds, $P(w) = 1/2$ (ajout/suppression)

Probabilité générique après k rounds (récurrence) : **$P(w) = 1/2^k$**

Nombre N de mots distincts du texte :

$$size(mem) \simeq N * P(w)$$

Un peu de probabilités



Probabilité $P(w)$ qu'un mot du texte soit dans la table après **2 rounds** (1 complet et 1 partiel) :

- si w n'apparaît qu'au premier round, $P(w) = 1/2$ (« nettoyage »)
- si w n'apparaît qu'au second round, $P(w) = 1/2$ (ajout)
- si w apparaît dans les 2 rounds, $P(w) = 1/2$ (ajout/suppression)

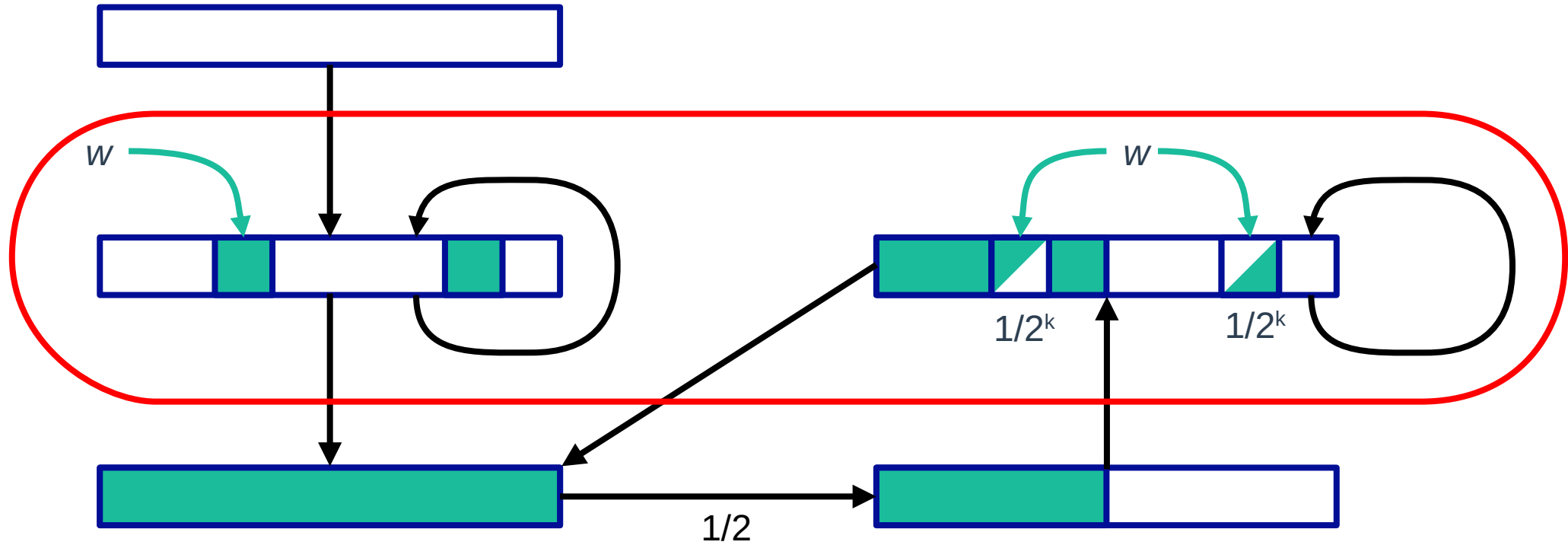
Probabilité générique après k rounds (réurrence) : **$P(w) = 1/2^k$**

Nombre N de mots distincts du texte :

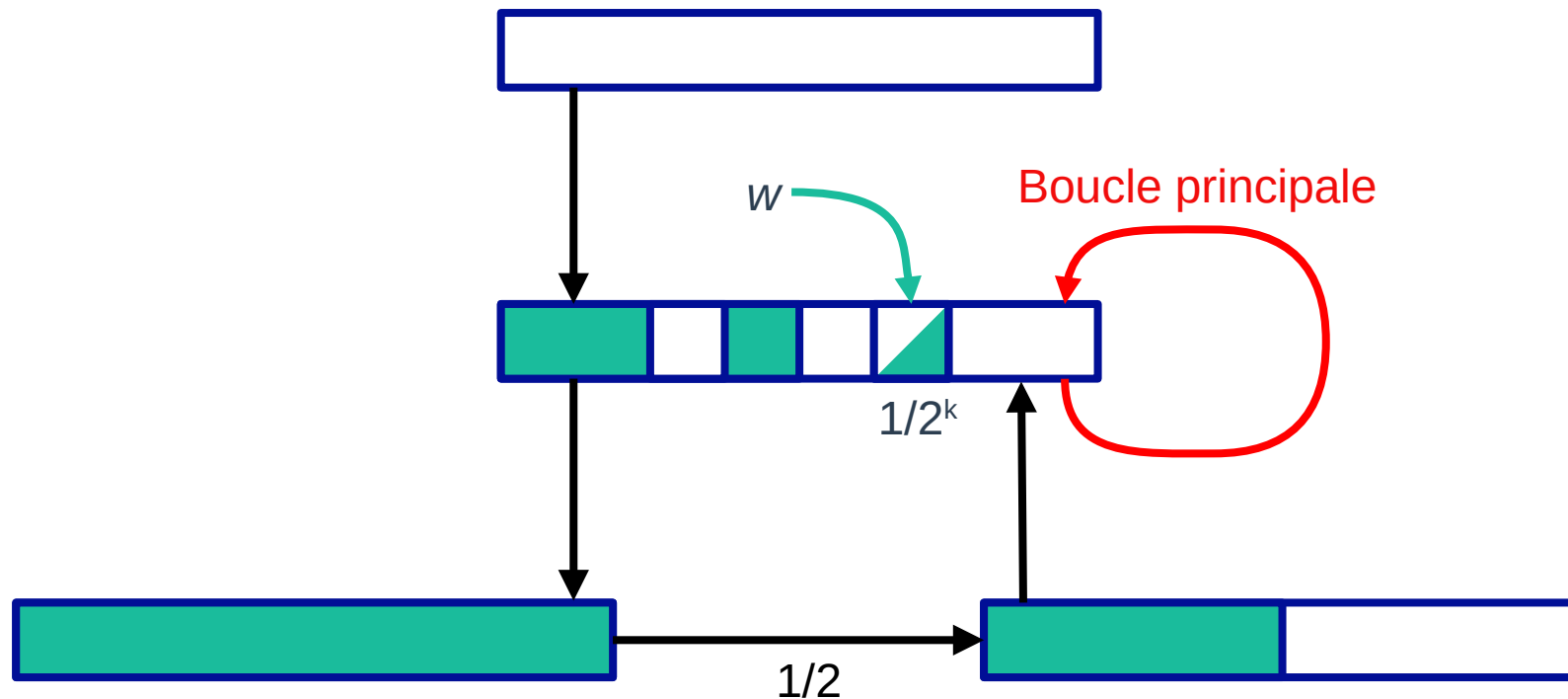
$$size(mem) \simeq N * P(w)$$

$$N \simeq size(mem) / P(w) = \mathbf{size(mem) * 2^k}$$

Implémentation



Implémentation



Décompte approximatif



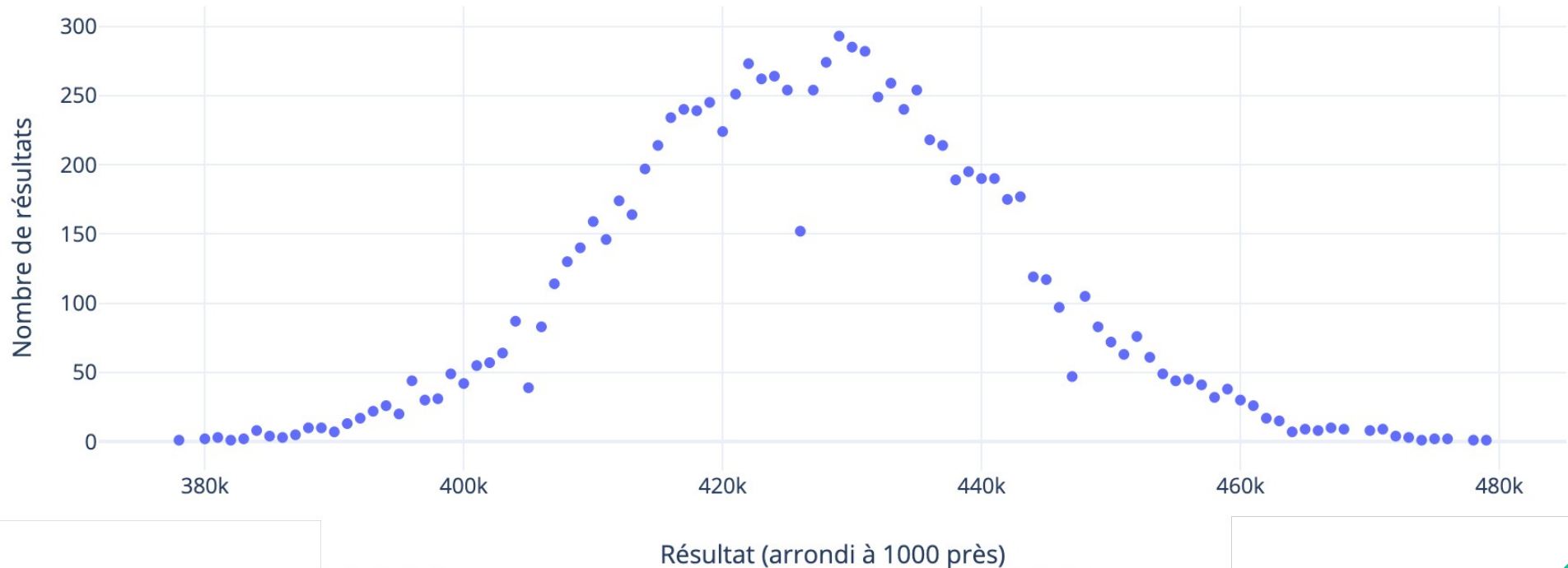
Décompte approximatif



- + performant
- + consommation mémoire maîtrisée
- = approximation raisonnable

Répartition des résultats

Répartition des résultats après 10 000 itérations (memSize = 1000)



Qualité de l'approximation (pour 427 430 mots distincts)

MemSize = 100

- **moyenne** = 427 717 (0.067% d'erreur)
- **écart-type** = 55 866 (13.06% de la moyenne)

66.7% des résultats ont moins de **13.06%** d'erreur

95% des résultats ont moins de 25.60% d'erreur

99.6% des résultats ont moins de 39.18% d'erreur

MemSize = 1 000

- **moyenne** = 427 576 (0.034% d'erreur)
- **écart-type** = 14 752 (3.45% de la moyenne)

66.7% des résultats ont moins de 3.45% d'erreur

95% des résultats ont moins de **6.76%** d'erreur

99.6% des résultats ont moins de **10.35%** d'erreur

MemSize = 10 000

- **moyenne** = 427 373 (0.013% d'erreur)
- **écart-type** = 5 109 (1.20% de la moyenne)

66.7% des résultats ont moins de 1.20% d'erreur

95% des résultats ont moins de 2.34% d'erreur

99.6% des résultats ont moins de **3.59%** d'erreur

Que peut-on faire d'autre ?



Pour aller plus loin...

Efficient probabilistic data structures (Raphael de Lio)

DEVOXX FRANCE 2025
3^{ème} édition - du 18 au 19 avril 2025



That's when I heard of Probabilistic Data Structures



DEVOXX FRANCE 2025
3^{ème} édition - du 18 au 19 avril 2025

PROBABILISTIC DATA STRUCTURES

- BLOOM FILTER
- COUNT-MIN SKETCH
- T-DIGEST
- HYPERLOGLOG
- MINHASH
- SIMHASH

Google Cloud | takima | THALES | W | AWS

@DEVOXXFR

Structures de données probabilistes (Hela Ben Khalfallah)

@DEVOXXFR



Google Cloud



Structures probabilistes

- Basées sur des **streams** et des **fonctions de hachage** (collisions)

Bloom Filter	Appartenance à un ensemble	<ul style="list-style-type: none">• mémoire bornée• recherche rapide• pas de faux négatifs
HyperLogLog	Nombre d'éléments uniques	<ul style="list-style-type: none">• mémoire limitée• combinable• peu d'erreurs ($\approx 1\%$)
Count-min Sketch	Nombre d'occurrences par élément	<ul style="list-style-type: none">• mémoire limitée• insertion rapide• surestime légèrement

source : *Hela Ben Khalfallah*

Probabilistic data structures

Probabilistic data structures let you answer common questions about your data streams and large datasets much faster. The significant advantages given in terms of efficiency in memory usage and processing speed is a trade-off with absolute accuracy. In Redis 8, in addition to **HyperLogLog**, which allows estimating the cardinality of a set, there are now 5 more probabilistic data structures:

- **Bloom** filter and **Cuckoo** filter—for checking if a given value has already appeared in a data stream
- **Count-min sketch**—for estimating how many times a given value appeared in a data stream
- **Top-k**—for finding the most frequent values in the data stream
- **t-digest**—for querying which fraction of the values in the data stream are smaller/larger than a given value.

Intégrations

- **Redis** : Bloom filter, Hyperloglog, Count-min sketch, Top-k, t-digest
- **PostgreSQL** : Bloom index
- **Spark** : Hyperloglog, Min-hash, Count-min sketch
- **Druid, Elasticsearch, BigQuery** : t-digest, Hyperloglog

source : Hela Ben Khalfallah

Conclusions



On peut encore faire de l'algo aujourd'hui
(et c'est pas forcément compliqué)



Les probabilités peuvent nous être utiles !



On n'a pas toujours besoin de résultats exacts
économisez des ressources si vous le pouvez



Merci !



Feedbacks



Code & Slides



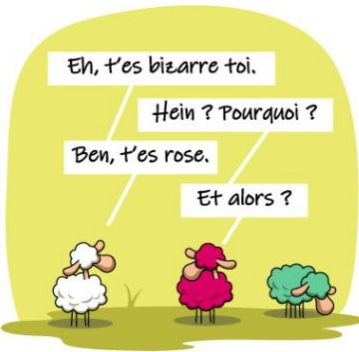
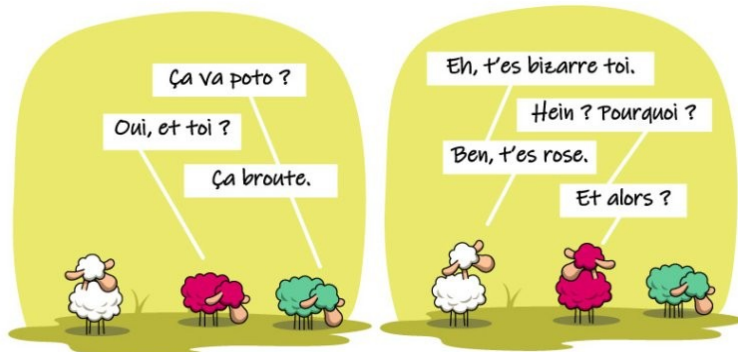
<https://github.com/benoitmasson/gi-go>



Feedbacks



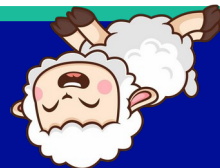
annyo.logaton.fr



Katam



(Bonus)



Code & Slides



<https://github.com/benoitmasson/gi-go>