

Spark et le streaming

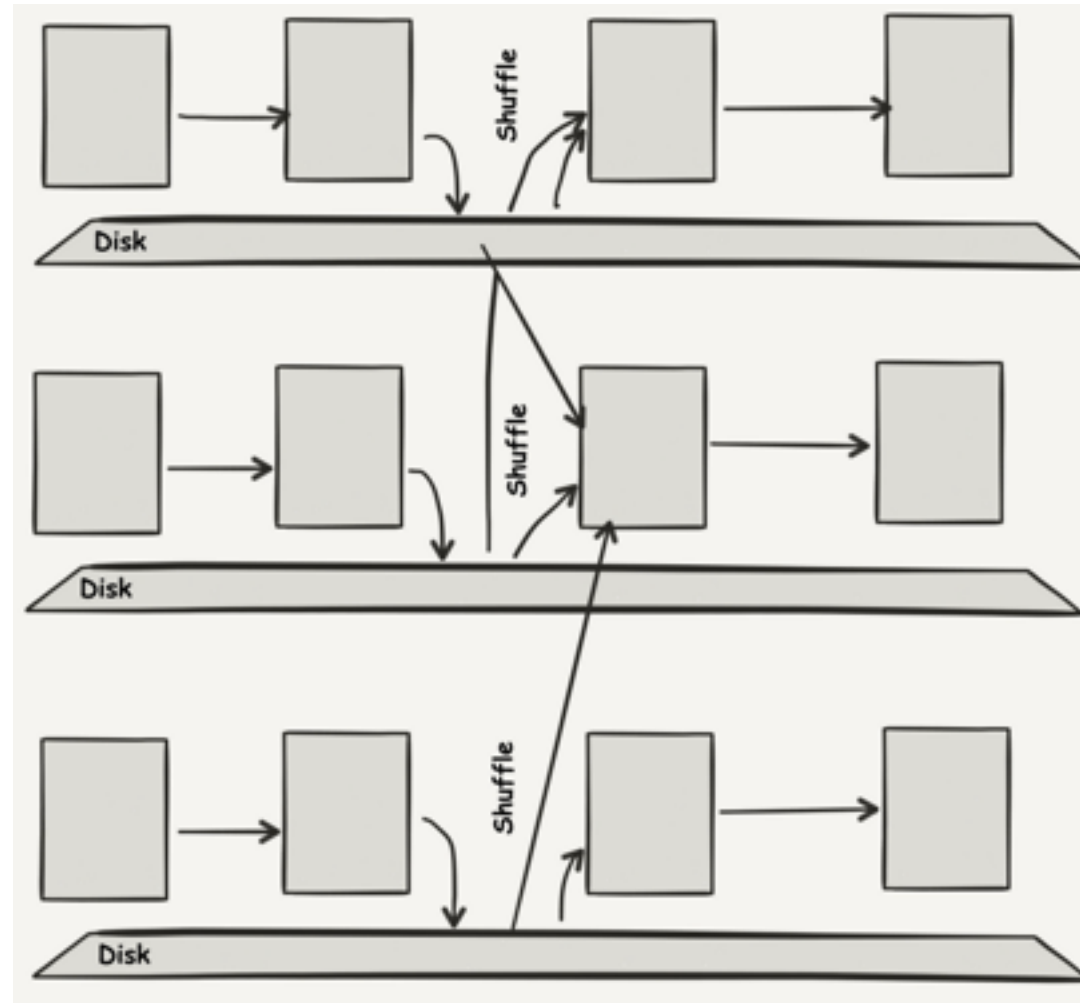
Jean-Paul LE
Geoffrey ALDEBERT

Optimisation

Optimisation

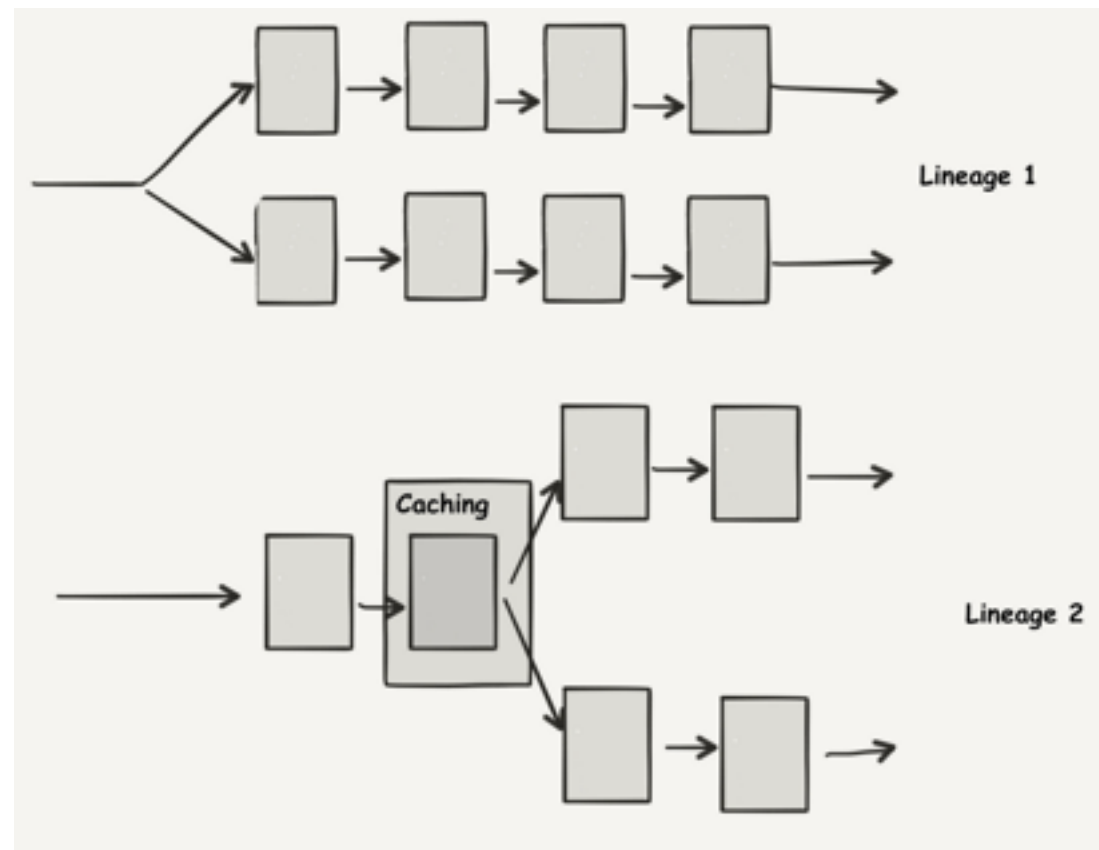
- Plan de calcul & Spark UI
- Partitionnement & Shuffling
- Caching & Checkpointing

Partitionnement & Shuffling



Caching & Checkpointing

- Caching/Persisting



Caching & Checkpointing

- Caching/Persisting
 - Ecrit en mémoire et/ou sur disque
 - Garde le lineage
 - Supprimé à la fermeture de l'application
- Type de sauvegarde
 - Sur disque
 - En mémoire
 - Sérialisé

Caching & Checkpointing

- Checkpointing
 - Ecrit sur disque
 - Supprime le lineage
 - Sauvegardé même après fermeture (utilisable par d'autre jobs)
 - Lent : mise en cache puis écriture sur disque recommandé

Caching & Checkpointing

- Sauvegarde de résultats intermédiaires pour réutilisation
- Plusieurs sauvegarde possibles :

	Caching	Checkpointing
Sauvegarde	Temporaire	Permanente
Cas d'utilisation	Réutilisation multiple d'un résultat	Sauvegarde d'un résultat après >100 transformations

Caching & Checkpointing

- TP

Spark Streaming

Spark Streaming

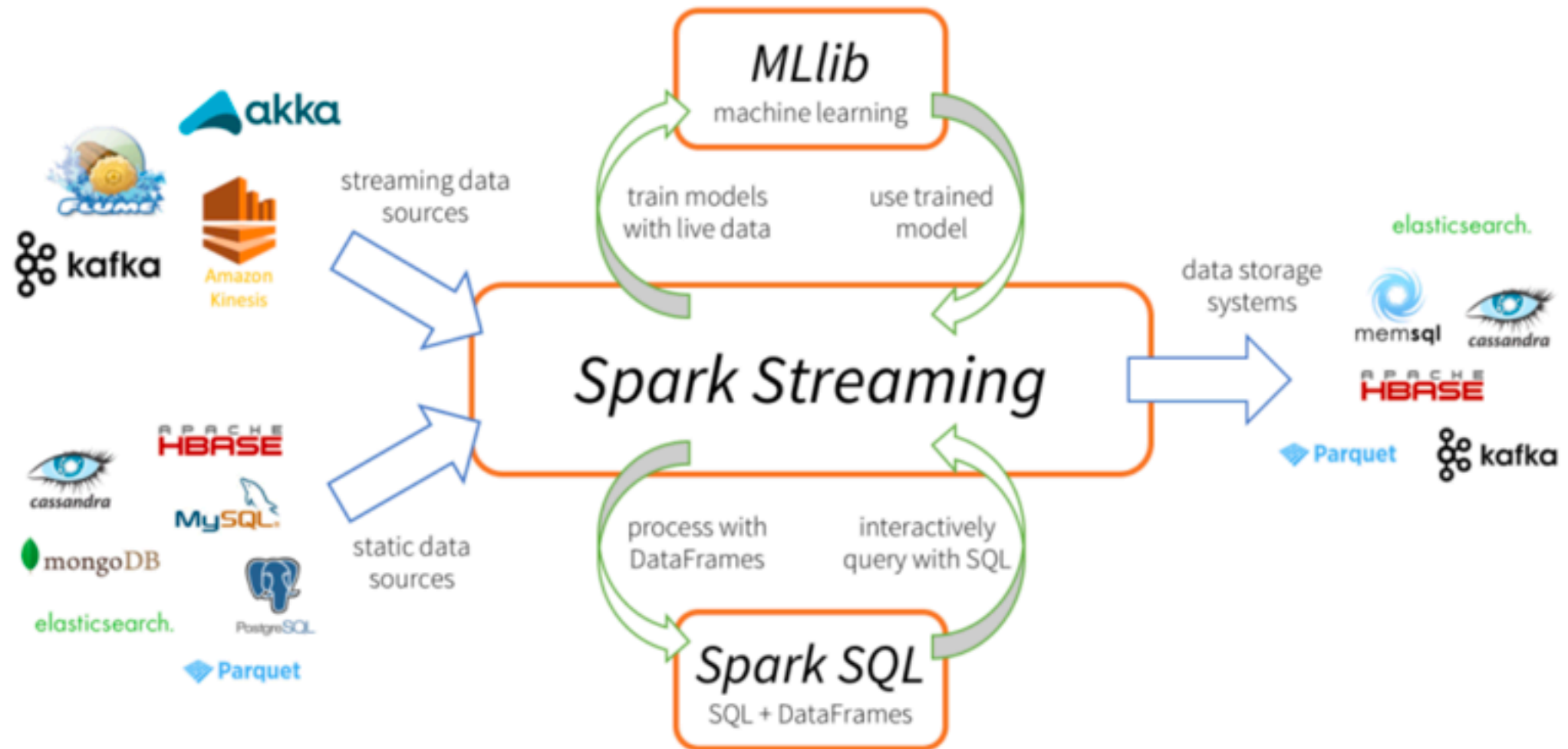


- Spark Streaming :
 - Récupère la donnée
 - Renvoie batch de données (DStream)
- Spark Engine
 - Manipulation des Stream

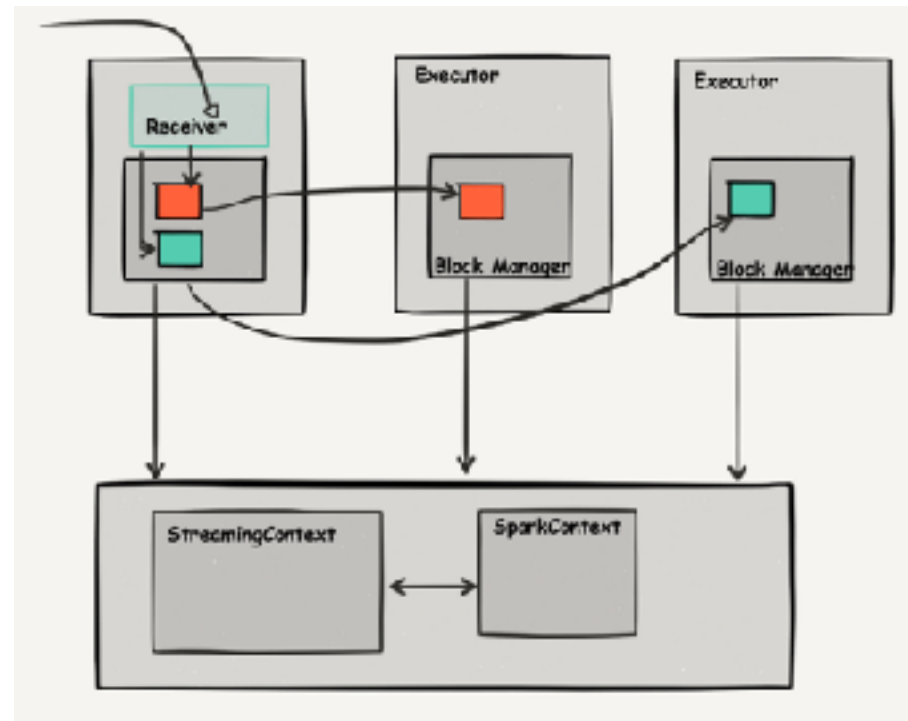
Spark Streaming

- Source de données:
 - Kafka / Flume / Kinesis / ...
 - HDFS / S3
 - TCP
 - Twitter

Spark Streaming

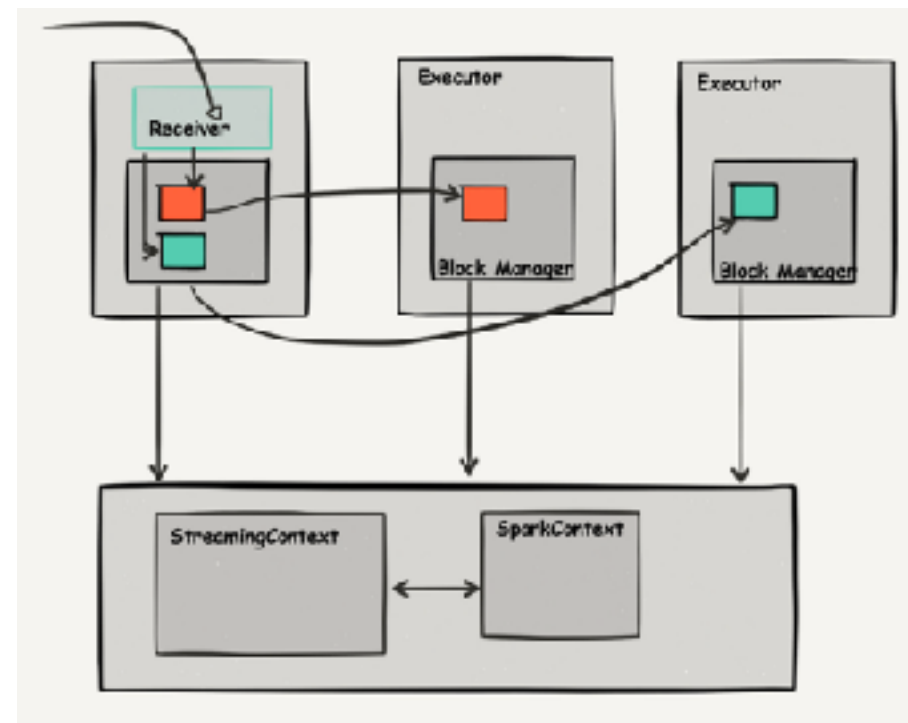


- Receivers

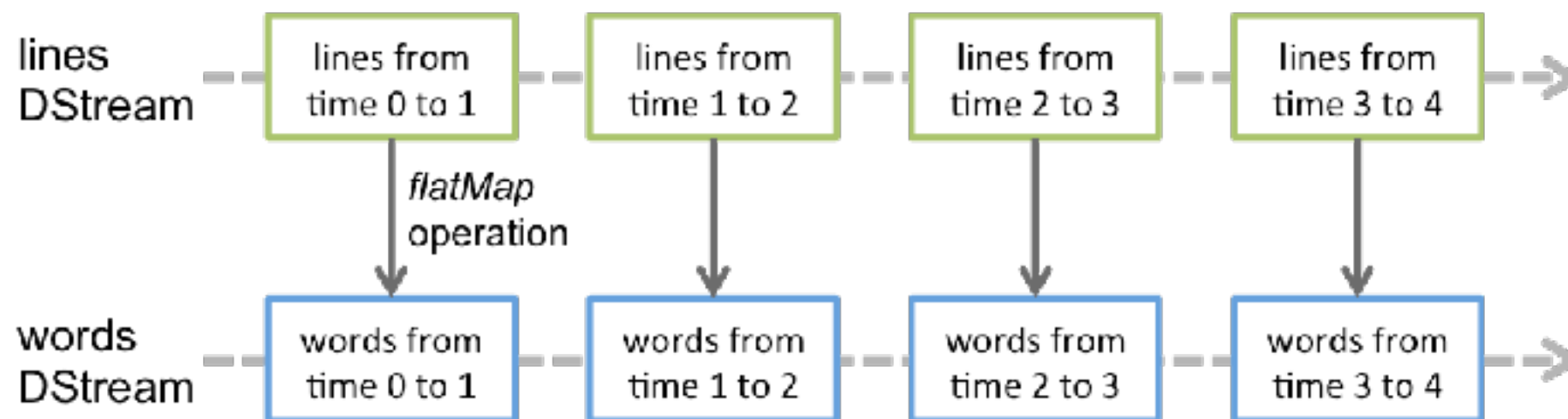


- DStream

- Receivers



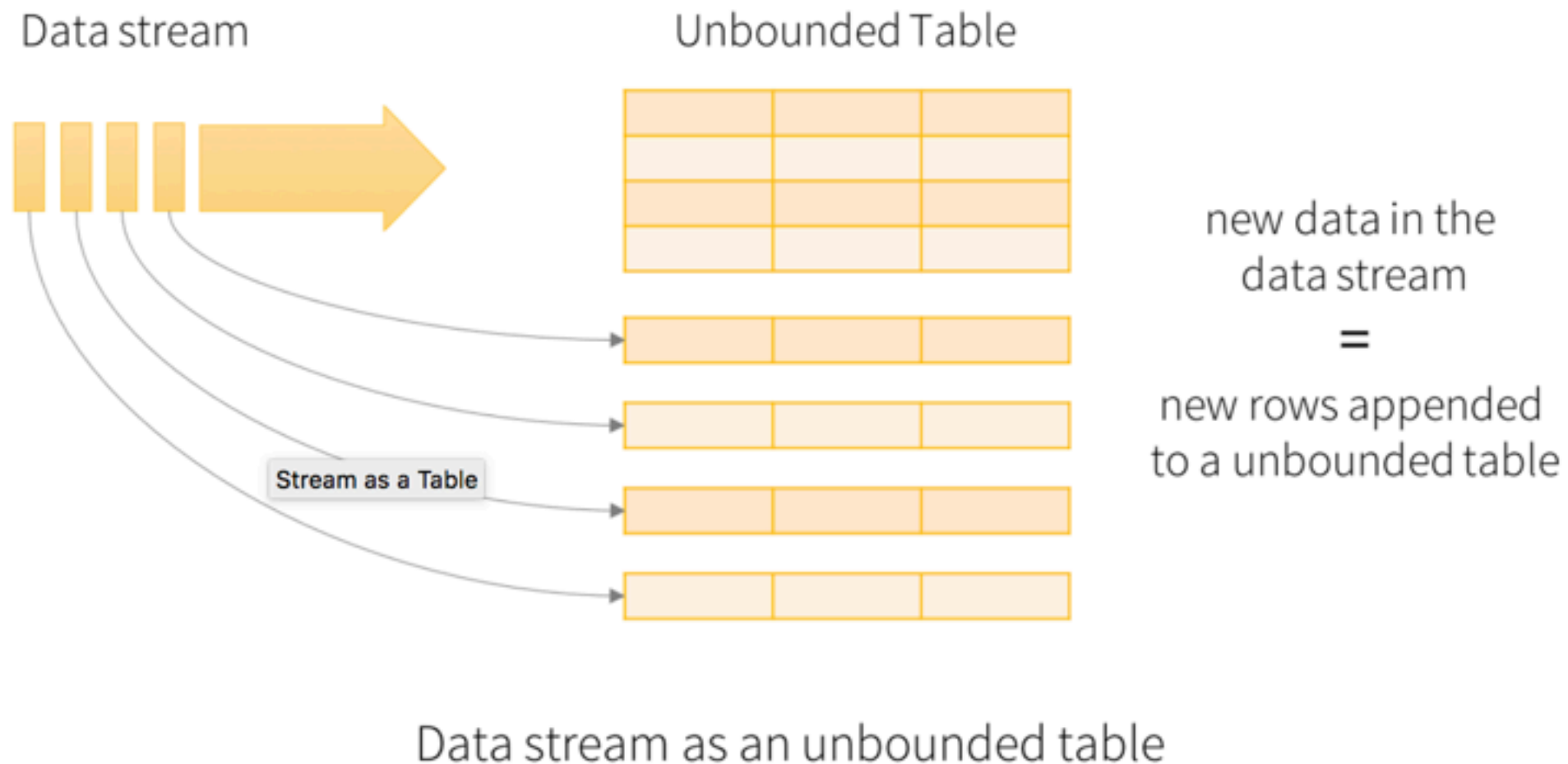
- DStream



Spark Structure Streaming

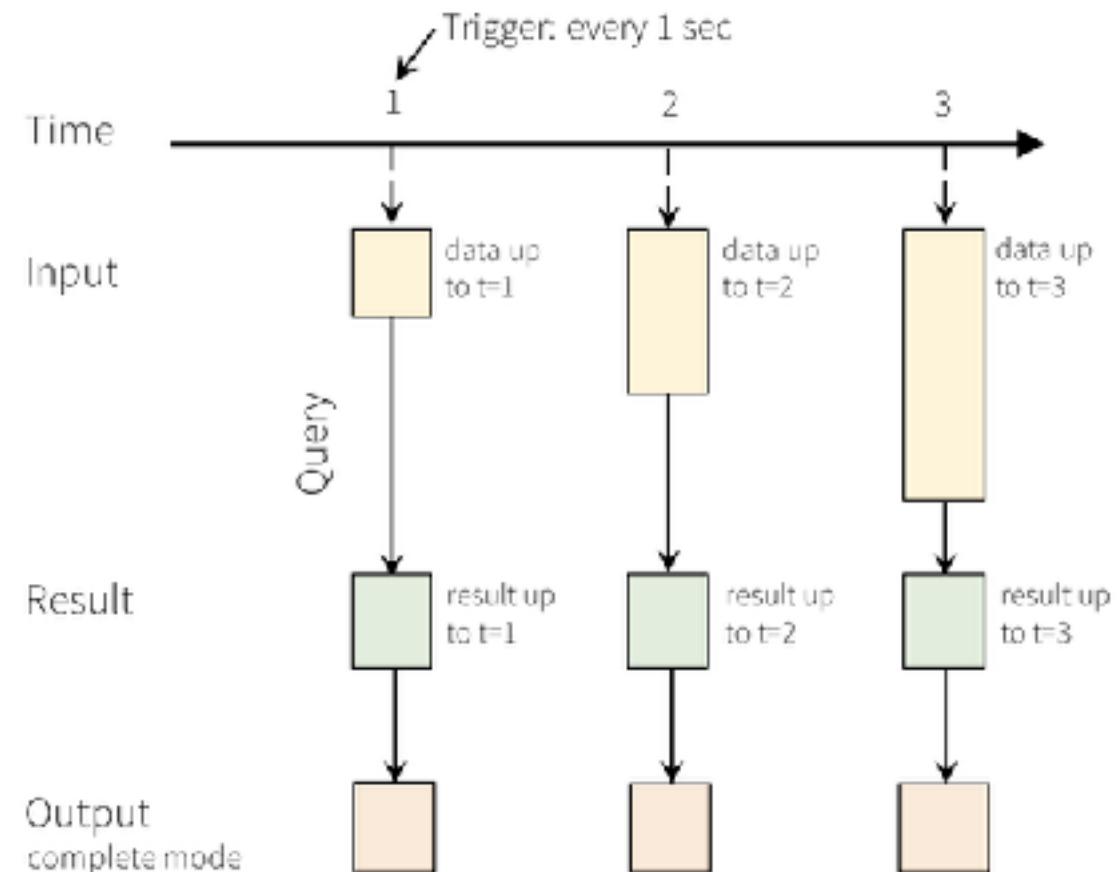
- RDD > Spark SQL (Dataframe)
- Spark Streaming > Spark Structured Streaming

Spark Structure Streaming



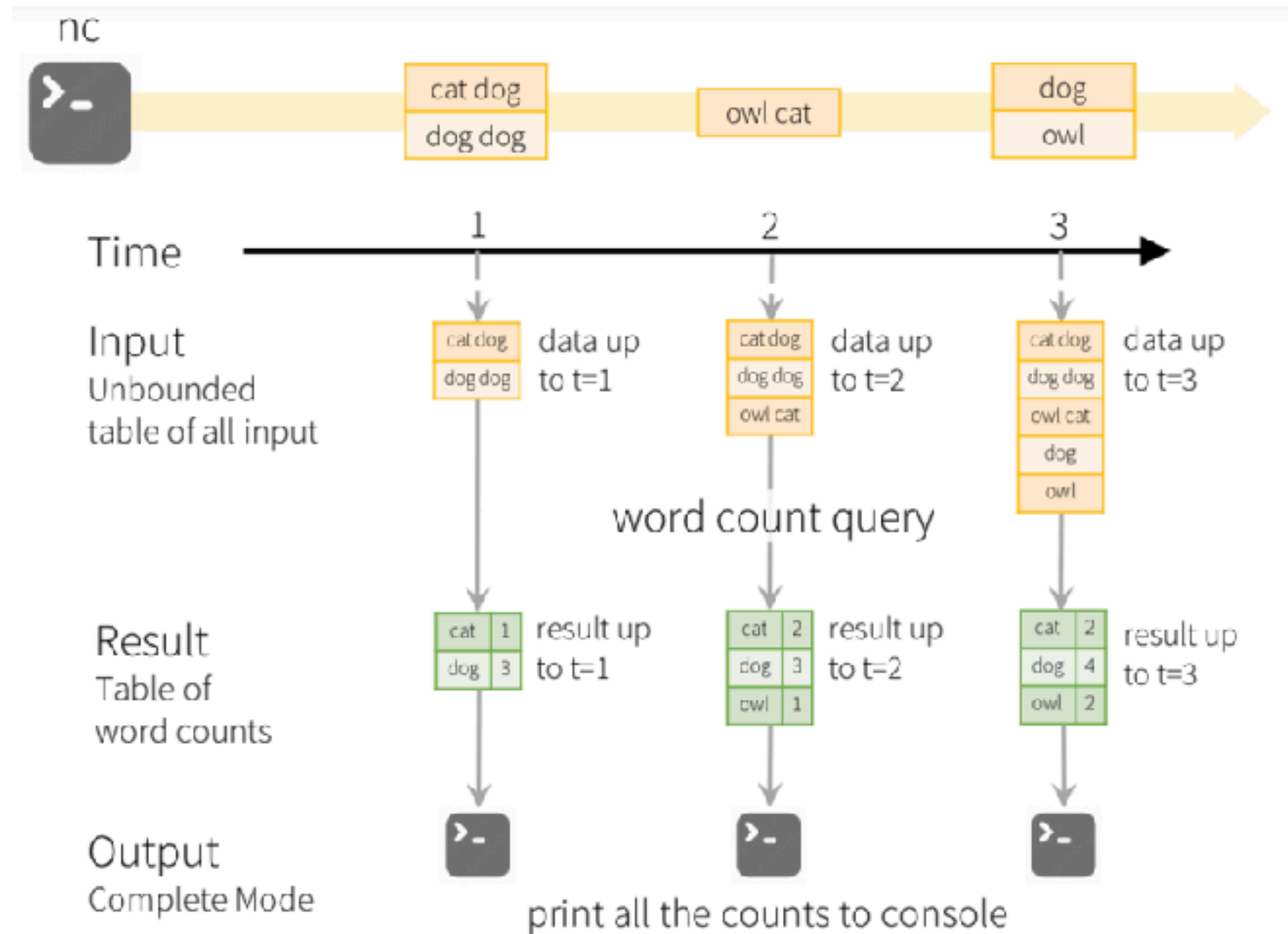
Spark Structure Streaming

- **Complete mode :**
Enregistre l'ensemble du tableau résultat
- **Append Mode :**
Enregistre que les nouvelles lignes
- **Update mode :**
Enregistre que les lignes mises à jour



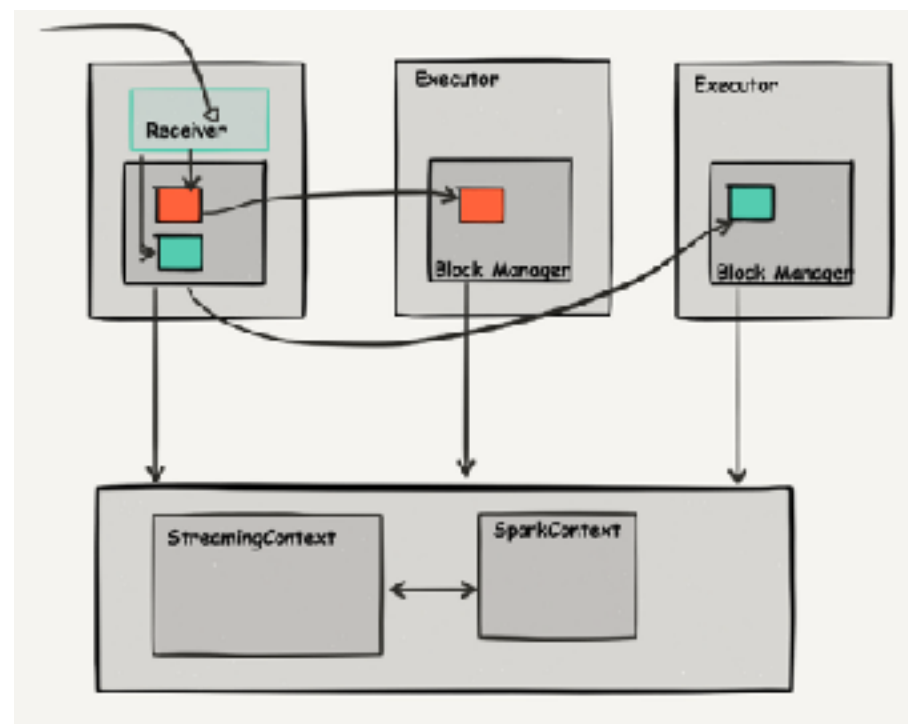
Programming Model for Structured Streaming

Spark Structure Streaming



Propriété de Spark Streaming

- Réplication des DStream
- Disponibilité du receiver



Intérêts de Spark Streaming

- Gère dynamiquement la charge des noeuds
- Tolérant à l'échec

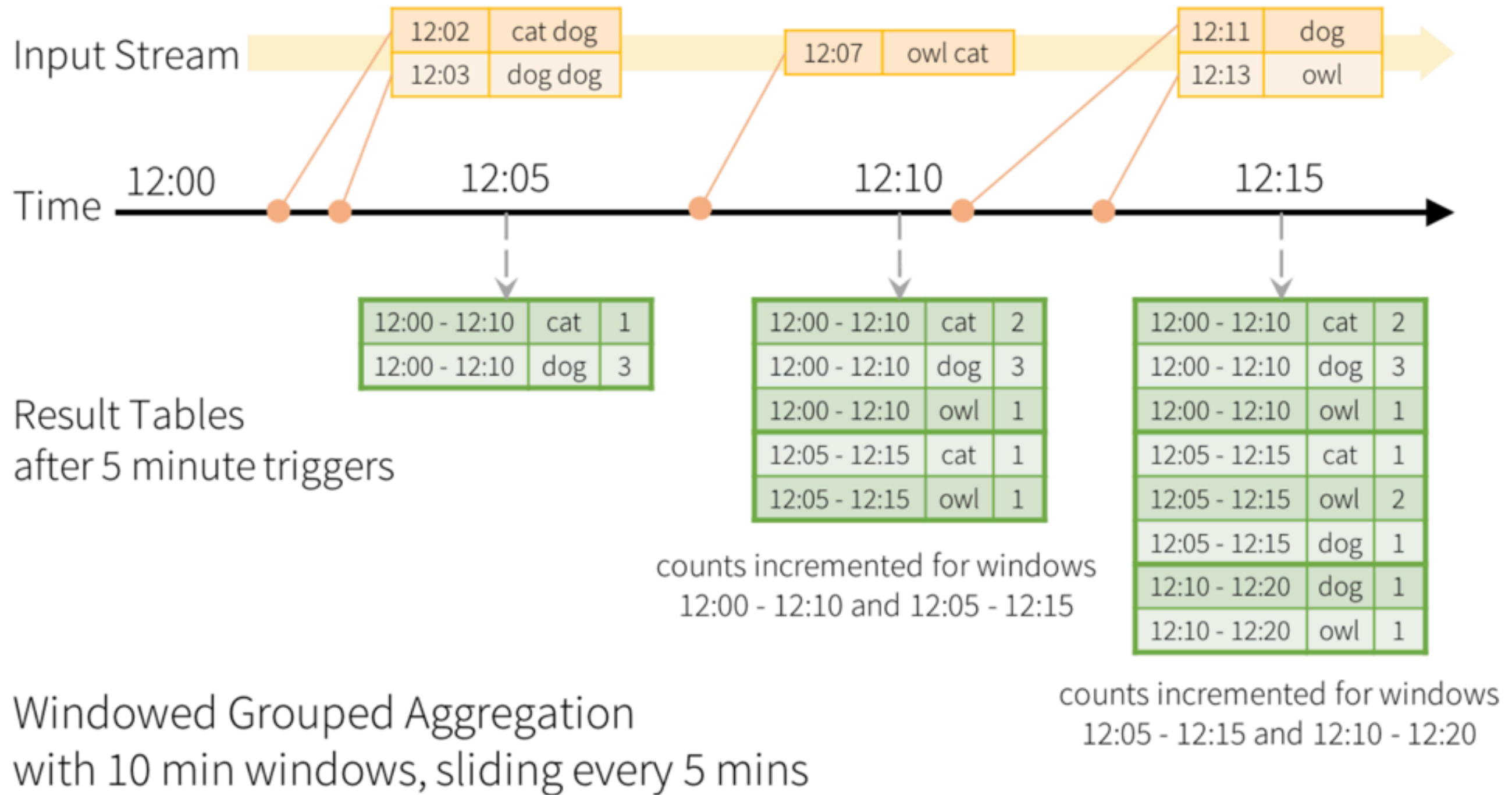
Limitation de Spark Streaming

- Receiver non fiable :
 - Pas de vérification de réception
 - Pas de récupération des données manquées
 - Perte de données possible
- Receiver custom en Scala/Java seulement

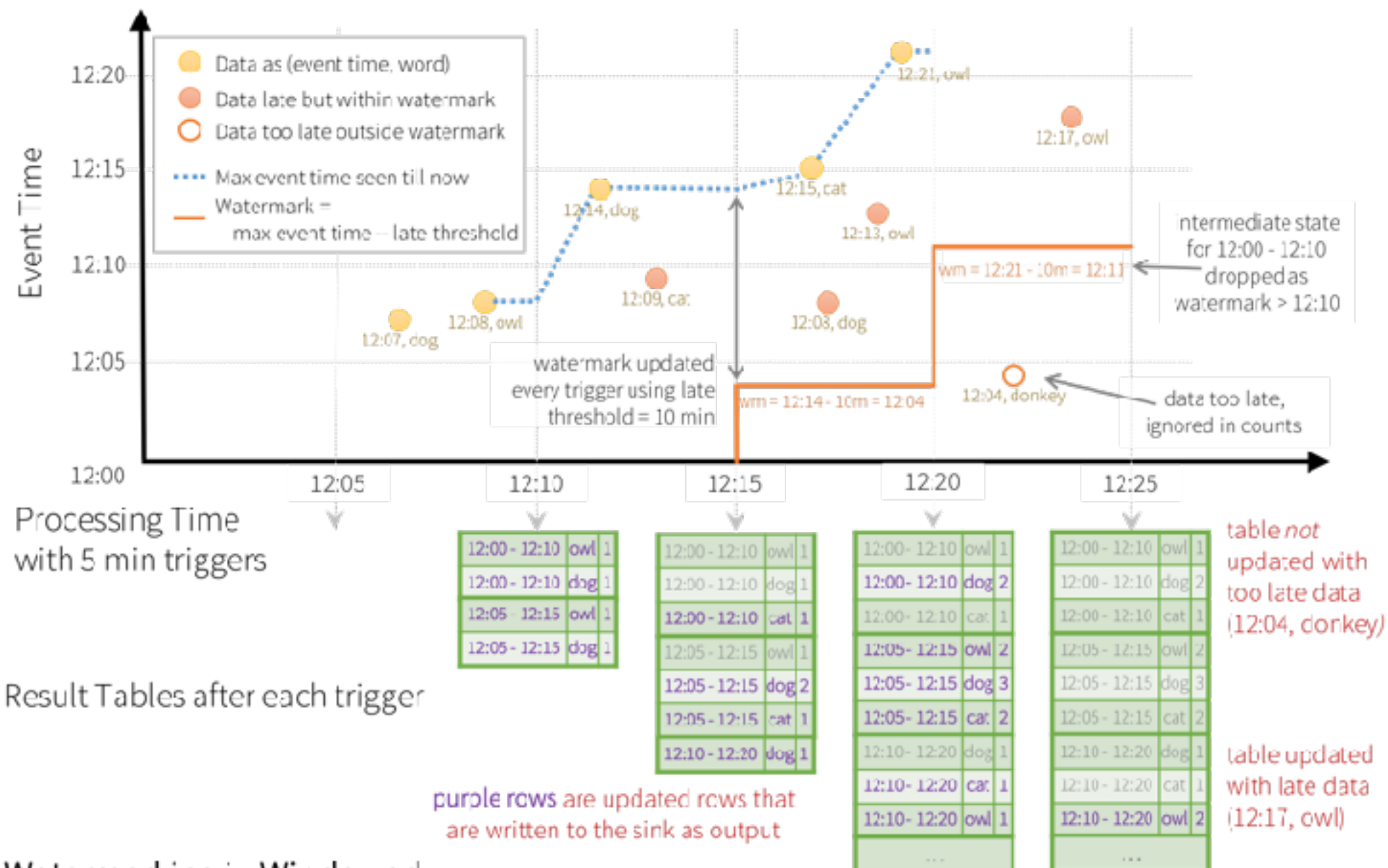
Spark Streaming

- TP

Fenêtrage



Fenêtrage (update mode)



Fenêtrage (append mode)

