

# Architecture Streaming Data

Architecture de stockage et de  
sortie des données en streaming

Amina MARIE

---







# Ordre du jour

Options de stockage pour les données en streaming

Mécanismes de sortie des données pour alimenter les systèmes en aval

Exercice 1 : Stockage des données en streaming dans un data lake à l'aide d'Apache Hadoop et Apache Parquet.

Exercice 2 : Visualisation des données en streaming à l'aide d'un outil de business intelligence tel que Tableau.



# Options de stockage pour les données en streaming



Stock Market Report

100.45	101.98	102.71	103.00	104.00
105.75	111.90	112.00	113.00	115.00
120.75	122.25	123.00	124.00	125.00
130.00	131.00	132.00	133.00	134.00
140.00	141.00	142.00	143.00	144.00
150.00	151.00	152.00	153.00	154.00
160.00	161.00	162.00	163.00	164.00
170.00	171.00	172.00	173.00	174.00
180.00	181.00	182.00	183.00	184.00
190.00	191.00	192.00	193.00	194.00
200.00	201.00	202.00	203.00	204.00
210.00	211.00	212.00	213.00	214.00
220.00	221.00	222.00	223.00	224.00
230.00	231.00	232.00	233.00	234.00
240.00	241.00	242.00	243.00	244.00
250.00	251.00	252.00	253.00	254.00
260.00	261.00	262.00	263.00	264.00
270.00	271.00	272.00	273.00	274.00
280.00	281.00	282.00	283.00	284.00
290.00	291.00	292.00	293.00	294.00
300.00	301.00	302.00	303.00	304.00
310.00	311.00	312.00	313.00	314.00
320.00	321.00	322.00	323.00	324.00
330.00	331.00	332.00	333.00	334.00
340.00	341.00	342.00	343.00	344.00
350.00	351.00	352.00	353.00	354.00
360.00	361.00	362.00	363.00	364.00
370.00	371.00	372.00	373.00	374.00
380.00	381.00	382.00	383.00	384.00
390.00	391.00	392.00	393.00	394.00
400.00	401.00	402.00	403.00	404.00
410.00	411.00	412.00	413.00	414.00
420.00	421.00	422.00	423.00	424.00
430.00	431.00	432.00	433.00	434.00
440.00	441.00	442.00	443.00	444.00
450.00	451.00	452.00	453.00	454.00
460.00	461.00	462.00	463.00	464.00
470.00	471.00	472.00	473.00	474.00
480.00	481.00	482.00	483.00	484.00
490.00	491.00	492.00	493.00	494.00
500.00	501.00	502.00	503.00	504.00
510.00	511.00	512.00	513.00	514.00
520.00	521.00	522.00	523.00	524.00
530.00	531.00	532.00	533.00	534.00
540.00	541.00	542.00	543.00	544.00
550.00	551.00	552.00	553.00	554.00
560.00	561.00	562.00	563.00	564.00
570.00	571.00	572.00	573.00	574.00
580.00	581.00	582.00	583.00	584.00
590.00	591.00	592.00	593.00	594.00
600.00	601.00	602.00	603.00	604.00
610.00	611.00	612.00	613.00	614.00
620.00	621.00	622.00	623.00	624.00
630.00	631.00	632.00	633.00	634.00
640.00	641.00	642.00	643.00	644.00
650.00	651.00	652.00	653.00	654.00
660.00	661.00	662.00	663.00	664.00
670.00	671.00	672.00	673.00	674.00
680.00	681.00	682.00	683.00	684.00
690.00	691.00	692.00	693.00	694.00
700.00	701.00	702.00	703.00	704.00
710.00	711.00	712.00	713.00	714.00
720.00	721.00	722.00	723.00	724.00
730.00	731.00	732.00	733.00	734.00
740.00	741.00	742.00	743.00	744.00
750.00	751.00	752.00	753.00	754.00
760.00	761.00	762.00	763.00	764.00
770.00	771.00	772.00	773.00	774.00
780.00	781.00	782.00	783.00	784.00
790.00	791.00	792.00	793.00	794.00
800.00	801.00	802.00	803.00	804.00
810.00	811.00	812.00	813.00	814.00
820.00	821.00	822.00	823.00	824.00
830.00	831.00	832.00	833.00	834.00
840.00	841.00	842.00	843.00	844.00
850.00	851.00	852.00	853.00	854.00
860.00	861.00	862.00	863.00	864.00
870.00	871.00	872.00	873.00	874.00
880.00	881.00	882.00	883.00	884.00
890.00	891.00	892.00	893.00	894.00
900.00	901.00	902.00	903.00	904.00
910.00	911.00	912.00	913.00	914.00
920.00	921.00	922.00	923.00	924.00
930.00	931.00	932.00	933.00	934.00
940.00	941.00	942.00	943.00	944.00
950.00	951.00	952.00	953.00	954.00
960.00	961.00	962.00	963.00	964.00
970.00	971.00	972.00	973.00	974.00
980.00	981.00	982.00	983.00	984.00
990.00	991.00	992.00	993.00	994.00
1000.00	1001.00	1002.00	1003.00	1004.00



# Options de stockage pour les données en streaming

Il n'y a pas de bonnes ou mauvaises options de stockage. Tout dépend des besoins du projet.

- Systèmes de fichiers distribués (DFS)
  - Bases de données NoSQL
  - Data Lakes
  - Bases de données In-Memory
  - Services de streaming gérés
-

# Options de stockage pour les données en streaming

- Systèmes de fichiers distribués (DFS)

Infrastructure informatique qui permet à plusieurs ordinateurs de travailler ensemble pour stocker et gérer des fichiers de manière distribuée sur un réseau.

Nous avons des DFS comme que Hadoop Distributed File System (HDFS) ou Amazon S3, qui sont conçus pour stocker de grandes quantités de données à l'échelle.

Ils offrent une redondance intégrée et une tolérance aux pannes, ce qui les rend robustes pour le stockage de données en streaming.

---

# Options de stockage pour les données en streaming

- Bases de données NoSQL

Les bases de données NoSQL, telles que MongoDB, Cassandra ou Apache CouchDB, sont bien adaptées au stockage de données en streaming en raison de leur capacité à gérer des flux de données volumineux avec une latence minimale.

Elles permettent également une évolutivité horizontale, ce qui signifie qu'elles peuvent s'adapter à une augmentation de la charge de travail.

---

# Options de stockage pour les données en streaming

- Data Lakes

Les Data Lakes, tels que Apache Hadoop ou Amazon S3, sont des dépôts de données qui stockent des données brutes dans leur format natif jusqu'à ce qu'elles soient nécessaires pour l'analyse.

Ils sont populaires pour le stockage de données en streaming car ils peuvent gérer une variété de données non structurées, semi-structurées et structurées.

---

# Options de stockage pour les données en streaming

- Bases de données In-Memory

Ces bases de données, telles que Redis ou Apache Ignite, stockent les données en mémoire vive plutôt que sur un disque, ce qui permet des temps de réponse très rapides.

Elles sont idéales pour les cas d'utilisation nécessitant un accès ultra-rapide aux données en streaming.

---



# Options de stockage pour les données en streaming

- Services de streaming gérés

Des services comme Amazon Kinesis ou Apache Kafka proposent des solutions intégrées pour la collecte, le traitement et le stockage des données en streaming.

Ils simplifient souvent la configuration et la gestion, mais peuvent être moins flexibles que les solutions auto-hébergées.

---

# Mécanismes de sortie des données pour alimenter les systèmes en aval



Stock Market Report

100.45	101.98	102.71	103.00	104.00
105.25	111.90	112.10	113.00	115.00
120.75	125.20	126.00	127.00	128.00
130.00	135.00	136.00	137.00	138.00
140.00	145.00	146.00	147.00	148.00
150.00	155.00	156.00	157.00	158.00
160.00	165.00	166.00	167.00	168.00
170.00	175.00	176.00	177.00	178.00
180.00	185.00	186.00	187.00	188.00
190.00	195.00	196.00	197.00	198.00
200.00	205.00	206.00	207.00	208.00
210.00	215.00	216.00	217.00	218.00
220.00	225.00	226.00	227.00	228.00
230.00	235.00	236.00	237.00	238.00
240.00	245.00	246.00	247.00	248.00
250.00	255.00	256.00	257.00	258.00
260.00	265.00	266.00	267.00	268.00
270.00	275.00	276.00	277.00	278.00
280.00	285.00	286.00	287.00	288.00
290.00	295.00	296.00	297.00	298.00
300.00	305.00	306.00	307.00	308.00
310.00	315.00	316.00	317.00	318.00
320.00	325.00	326.00	327.00	328.00
330.00	335.00	336.00	337.00	338.00
340.00	345.00	346.00	347.00	348.00
350.00	355.00	356.00	357.00	358.00
360.00	365.00	366.00	367.00	368.00
370.00	375.00	376.00	377.00	378.00
380.00	385.00	386.00	387.00	388.00
390.00	395.00	396.00	397.00	398.00
400.00	405.00	406.00	407.00	408.00
410.00	415.00	416.00	417.00	418.00
420.00	425.00	426.00	427.00	428.00
430.00	435.00	436.00	437.00	438.00
440.00	445.00	446.00	447.00	448.00
450.00	455.00	456.00	457.00	458.00
460.00	465.00	466.00	467.00	468.00
470.00	475.00	476.00	477.00	478.00
480.00	485.00	486.00	487.00	488.00
490.00	495.00	496.00	497.00	498.00
500.00	505.00	506.00	507.00	508.00
510.00	515.00	516.00	517.00	518.00
520.00	525.00	526.00	527.00	528.00
530.00	535.00	536.00	537.00	538.00
540.00	545.00	546.00	547.00	548.00
550.00	555.00	556.00	557.00	558.00
560.00	565.00	566.00	567.00	568.00
570.00	575.00	576.00	577.00	578.00
580.00	585.00	586.00	587.00	588.00
590.00	595.00	596.00	597.00	598.00
600.00	605.00	606.00	607.00	608.00
610.00	615.00	616.00	617.00	618.00
620.00	625.00	626.00	627.00	628.00
630.00	635.00	636.00	637.00	638.00
640.00	645.00	646.00	647.00	648.00
650.00	655.00	656.00	657.00	658.00
660.00	665.00	666.00	667.00	668.00
670.00	675.00	676.00	677.00	678.00
680.00	685.00	686.00	687.00	688.00
690.00	695.00	696.00	697.00	698.00
700.00	705.00	706.00	707.00	708.00
710.00	715.00	716.00	717.00	718.00
720.00	725.00	726.00	727.00	728.00
730.00	735.00	736.00	737.00	738.00
740.00	745.00	746.00	747.00	748.00
750.00	755.00	756.00	757.00	758.00
760.00	765.00	766.00	767.00	768.00
770.00	775.00	776.00	777.00	778.00
780.00	785.00	786.00	787.00	788.00
790.00	795.00	796.00	797.00	798.00
800.00	805.00	806.00	807.00	808.00
810.00	815.00	816.00	817.00	818.00
820.00	825.00	826.00	827.00	828.00
830.00	835.00	836.00	837.00	838.00
840.00	845.00	846.00	847.00	848.00
850.00	855.00	856.00	857.00	858.00
860.00	865.00	866.00	867.00	868.00
870.00	875.00	876.00	877.00	878.00
880.00	885.00	886.00	887.00	888.00
890.00	895.00	896.00	897.00	898.00
900.00	905.00	906.00	907.00	908.00
910.00	915.00	916.00	917.00	918.00
920.00	925.00	926.00	927.00	928.00
930.00	935.00	936.00	937.00	938.00
940.00	945.00	946.00	947.00	948.00
950.00	955.00	956.00	957.00	958.00
960.00	965.00	966.00	967.00	968.00
970.00	975.00	976.00	977.00	978.00
980.00	985.00	986.00	987.00	988.00
990.00	995.00	996.00	997.00	998.00
1000.00	1005.00	1006.00	1007.00	1008.00

# Mécanismes de sortie des données pour alimenter les systèmes en aval

- Quels sont les systèmes en aval?
  - Visualisation
  - Tableaux de bord
  - Base de données...

Les mécanismes:

- Kafka Connect
  - API REST
  - Intégrations natives
  - Protocoles de streaming
  - Connecteurs personnalisés
-



# Mécanismes de sortie des données pour alimenter les systèmes en aval

- Kafka Connect

Kafka Connect est une plateforme open source pour connecter Kafka à des systèmes externes tels que des bases de données, des entrepôts de données, etc. Il permet d'importer et d'exporter des données en temps réel vers et depuis Kafka.

- API REST

De nombreuses plateformes de streaming offrent des API REST pour permettre aux utilisateurs d'envoyer et de recevoir des données en temps réel. Ces API sont souvent utilisées pour intégrer des applications tierces avec des systèmes de streaming.

---

# Mécanismes de sortie des données pour alimenter les systèmes en aval

- Intégrations natives

Certains systèmes de streaming ont des intégrations natives avec d'autres technologies courantes telles que Hadoop, Spark, Elasticsearch, etc. Ces intégrations facilitent le transfert de données entre différents systèmes sans nécessiter de développements supplémentaires.

- Protocoles de streaming

Des protocoles comme MQTT, AMQP, et d'autres sont souvent utilisés pour le streaming de données en temps réel. Ils permettent de transférer des données de manière efficace et fiable entre les différents composants du système.

---

# Mécanismes de sortie des données pour alimenter les systèmes en aval

- Connecteurs personnalisés

Dans certains cas, des connecteurs personnalisés peuvent être développés pour intégrer des systèmes spécifiques avec des plateformes de streaming. Ces connecteurs sont généralement développés en fonction des besoins spécifiques de l'application.

L'architecture de sortie des données en streaming doit être conçue avec soin pour garantir la fiabilité, la scalabilité et les performances nécessaires pour alimenter efficacement les systèmes en aval.

---



Stock Market Report

Stock	Price	Change	Volume
AAPL	141.04	+0.71	142,981
MSFT	105.25	+0.18	108,001
GOOGL	102.73	+0.15	102,001
AMZN	101.04	+0.12	101,001
FB	100.00	+0.10	100,001
BRK.A	100.00	+0.10	100,001
WMT	100.00	+0.10	100,001
DIS	100.00	+0.10	100,001
IBM	100.00	+0.10	100,001
ORCL	100.00	+0.10	100,001
CRM	100.00	+0.10	100,001
ADBE	100.00	+0.10	100,001
QCOM	100.00	+0.10	100,001
TXN	100.00	+0.10	100,001
INTC	100.00	+0.10	100,001
AMD	100.00	+0.10	100,001
AVGO	100.00	+0.10	100,001
MRNA	100.00	+0.10	100,001
REGN	100.00	+0.10	100,001
VRTX	100.00	+0.10	100,001
BIIB	100.00	+0.10	100,001
CRMR	100.00	+0.10	100,001
MRK	100.00	+0.10	100,001
ABBV	100.00	+0.10	100,001
LLY	100.00	+0.10	100,001
UNH	100.00	+0.10	100,001
CVS	100.00	+0.10	100,001
ANTM	100.00	+0.10	100,001
COO	100.00	+0.10	100,001
ABT	100.00	+0.10	100,001
MRN	100.00	+0.10	100,001
MDT	100.00	+0.10	100,001
BDX	100.00	+0.10	100,001
BSX	100.00	+0.10	100,001
EW	100.00	+0.10	100,001
OMC	100.00	+0.10	100,001
WDC	100.00	+0.10	100,001
SPGI	100.00	+0.10	100,001
MO	100.00	+0.10	100,001
CL	100.00	+0.10	100,001
PG	100.00	+0.10	100,001
KO	100.00	+0.10	100,001
PEP	100.00	+0.10	100,001
DIS	100.00	+0.10	100,001
CMSCA	100.00	+0.10	100,001
CSX	100.00	+0.10	100,001
CS	100.00	+0.10	100,001
NSR	100.00	+0.10	100,001
OTIS	100.00	+0.10	100,001
CTSH	100.00	+0.10	100,001
CTA	100.00	+0.10	100,001
CTI	100.00	+0.10	100,001
CTO	100.00	+0.10	100,001
CTP	100.00	+0.10	100,001
CTQ	100.00	+0.10	100,001
CTR	100.00	+0.10	100,001
CTS	100.00	+0.10	100,001
CTT	100.00	+0.10	100,001
CTV	100.00	+0.10	100,001
CTW	100.00	+0.10	100,001
CTX	100.00	+0.10	100,001
CTY	100.00	+0.10	100,001
CTZ	100.00	+0.10	100,001
CTAA	100.00	+0.10	100,001
CTAB	100.00	+0.10	100,001
CTAC	100.00	+0.10	100,001
CTAD	100.00	+0.10	100,001
CTAE	100.00	+0.10	100,001
CTAF	100.00	+0.10	100,001
CTAG	100.00	+0.10	100,001
CTAH	100.00	+0.10	100,001
CTAI	100.00	+0.10	100,001
CTAJ	100.00	+0.10	100,001
CTAK	100.00	+0.10	100,001
CTAL	100.00	+0.10	100,001
CTAM	100.00	+0.10	100,001
CTAN	100.00	+0.10	100,001
CTAO	100.00	+0.10	100,001
CTAP	100.00		

# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

Tableau est un outil de business intelligence (BI) très populaire.

Il permet aux utilisateurs de visualiser et d'analyser leurs données de manière interactive.

**Visualisation des données :** tableaux de bord interactifs, des graphiques, des cartes, des tableaux croisés dynamiques..

**Connexion à diverses sources de données :** il peut se connecter à une variété de sources de données (bases de données relationnelles, NoSQL, des fichiers plats, services cloud, des données en streaming, etc.)

**Analyse avancée :** Tableau offre des fonctionnalités avancées d'analyse des données, telles que la prévision, le clustering, les calculs personnalisés, l'agrégation de données, etc., permettant aux utilisateurs d'extraire des insights approfondis à partir de leurs données.

---

# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

**Interactivité :** Les utilisateurs peuvent explorer les données de manière interactive en filtrant, en forant, en zoomant et en sélectionnant des parties spécifiques des visualisations. Cela permet une analyse approfondie des données sans avoir besoin de compétences techniques avancées en requêtes SQL ou en programmation.

**Tableaux de bord dynamiques :** Tableau permet de créer des tableaux de bord dynamiques qui regroupent plusieurs visualisations pour offrir une vue d'ensemble complète des données. Les utilisateurs peuvent interagir avec les visualisations et les filtres pour obtenir des insights en temps réel.

**Partage et collaboration :** Les tableaux de bord et les visualisations créés dans Tableau peuvent être partagés avec d'autres utilisateurs via des liens web, des fichiers téléchargeables ou des intégrations dans des plates-formes de collaboration telles que Slack ou Microsoft Teams.

---



# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

Objectif: Utiliser Tableau et créer un tableau de visualisation de données.

Les étapes:

- Aller sur le site de Tableau <https://www.tableau.com/> et créer un compte d'essai.
- Activer votre compte et aller sur votre interface utilisateur.
- Créer un flux de données streaming qui représentant les ventes d'une entreprise fictive (coté producteur).  
Utilisez Kafka.

Les données doivent inclure: Date/heure de la vente, Montant de la vente, Produit vendu, Région de vente

- Sauvegarder les données cotés consommateur dans une base de données NoSQL. Utilisez MongoDB
-

# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

- Créer un projet Tableau et aller dans ce projet
- Configurer votre flux de données à Tableau en lui indiquant le serveur et le port sur lequel il doit écouter.
- Une fois connecté à vos données en streaming, commencez à créer vos visualisations.

Utilisez les champs de données appropriés pour créer des graphiques représentant l'évolution du chiffre d'affaires, la répartition géographique des ventes, les produits les plus vendus, etc.

- Explorez les différentes options de visualisation disponibles dans Tableau pour choisir les meilleures représentations pour vos données.
-

# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

- Ajout d'interactions et de filtres : Ajoutez des filtres interactifs pour permettre aux utilisateurs de sélectionner une période de temps spécifique, une région spécifique ou un produit spécifique.
  - Personnalisation du tableau de bord : Créez un tableau de bord dans Tableau Desktop en rassemblant toutes vos visualisations créées précédemment.
-



# Exercice 2 : Visualisation des données en streaming à l'aide de Tableau.

Ce qui est attendu:

- Un dépôt git contenant le producteur et le consommateur.
- Un compte rendu :
  - Expliquant la démarche mise en œuvre durant l'exercice
  - Contenant des screens du projet Tableau
  - Les différentes visualisations produites.
- Un partage du projet Tableau avec moi-même pour que je puisse y accéder.
- Une mise à jour en temps réel correcte des visualisations par rapport aux données.
- Des filtres interactifs fonctionnel comme prévu et des visualisations qui reflètent les sélections des utilisateurs.

Cet exercice compte pour l'évaluation finale. A rendre en binôme avant 21h. Aucun projet ne sera accepté passé ce délai.

---