



# ARQUITECTURAS BIG DATA, APACHE SPARK Y KAFKA

## Actividad de transferencia

### Definiendo Tecnologías

En los proyectos de Big Data y fundamentalmente en aquellos proyectos que son de Big Data Streaming, es decir, de tiempo casi-real, hemos visto que existen arquitecturas por un lado y tecnologías por otro que permiten de alguna forma dar respuesta eficiente a estos proyectos.

Entre las arquitecturas vimos la arquitectura Lambda y la Arquitectura Kappa. Entre las tecnologías, vimos Apache Spark y Apache Kafka. No obstante, es frecuente que en los proyectos de big data se combinen diversas tecnologías.

La actividad de esta semana consiste en armar un **cuadro comparativo**, para que puedas conocer otras tecnologías también muy utilizadas y, sobre todo, para revisar algunas diferencias sutiles que hay entre ellas.

Vamos a comparar las siguientes tecnologías:

- Apache Storm
- Apache Spark
- Apache Samza

Aplicación	¿Para qué se utiliza?	Algunas ventajas	Proyectos ideales para implementar esta tecnología	Ejemplo de empresa que lo utiliza
Apache Storm	Procesamiento de datos en tiempo real y en flujos continuos.	<ul style="list-style-type: none"><li>• Baja latencia.</li><li>• Alto rendimiento</li></ul>	<ul style="list-style-type: none"><li>• Análisis de redes sociales en</li></ul>	Twitter



		<ul style="list-style-type: none"><li>• Escalabilidad horizontal</li><li>• Arquitectura basada en topologías flexibles</li></ul>	<p>tiempo real</p> <ul style="list-style-type: none"><li>• Monitoreo de sistemas y detección de fraudes</li><li>• Procesamiento de logs y clics</li></ul>	
<b>Apache Spark</b>	Procesamiento en tiempo real (con Spark Streaming), por lotes y análisis avanzado de grandes volúmenes de datos.	<ul style="list-style-type: none"><li>• Soporte para múltiples lenguajes (java, scala, python y R)</li><li>• Alto rendimiento con procesamiento en memoria</li><li>• API unificada para batch y streaming</li></ul>	<ul style="list-style-type: none"><li>• Machine learning con grandes volúmenes de datos</li><li>• Procesamiento ETL</li><li>• Análisis de logs, métricas y comportamiento de usuarios</li></ul>	Netflix, eBay
<b>Apache Samza</b>	Procesamiento de flujos en tiempo real, especialmente con integración nativa a Kafka y YARN.	<ul style="list-style-type: none"><li>• Integración nativa con Apache Kafka</li><li>• Alta tolerancia a fallos</li><li>• Procesamiento local</li></ul>	<ul style="list-style-type: none"><li>• Aplicaciones de monitoreo en tiempo real</li><li>• Dashboards de</li></ul>	LinkedIn



		(near data processing)	eventos en vivo <ul style="list-style-type: none"><li>● Enriquecimiento de datos en flujos Kafka</li></ul>	
--	--	------------------------	--	--

#### Fuente de datos consultada:

- Apache Storm Official Documentation: <https://storm.apache.org/>
- Apache Spark Official Documentation: <https://spark.apache.org/>
- Apache Samza Official Documentation: <https://samza.apache.org/>
- Karau, H., & Warren, R. (2015). *High Performance Spark*. O'Reilly Media.
- Gulisano, V., et al. (2012). *StreamCloud: An elastic and scalable data streaming system*. IEEE Transactions on Parallel and Distributed Systems.
- Zaharia, M. et al. (2016). *Apache Spark: a unified engine for big data processing*. Communications of the ACM.
- Toshniwal, A. et al. (2014). *Storm@Twitter: Scaling distributed stream processing*. VLDB.
- Noghabi, S.A. et al. (2017). *Samza: stateful scalable stream processing at LinkedIn*. Proceedings of the VLDB Endowment.
- Netflix Tech Blog: <https://netflixtechblog.com/>

La entrega será un documento WORD o PDF con el cuadro comparativo. Debes de estar lo más completo que puedas y, sobre todo, debes colocar las referencias de dónde has obtenido la información.

NOTA: Recuerda utilizar recursos como Google Académico, para acceder a papers o documentos académicos que puedan tener validez para educación.