

Red Hat  
**Summit**

## Connect

# Integración de tareas de calidad de datos con IA Generativa en MLOps

Marlon Cárdenas, PhD.

Área de Data & IA

sopra  steria

# Agenda

- La Era de “Conversar” con los datos
- **Apuesta de la compañía:** Plataforma Inteligente
- **Nuevo objetivo:** usar, entrenar y evaluar LLMs
- **Nuevos retos:** generar datos sintéticos en contextos sensibles
- **Discusión:** problemas con el uso de datos sintéticos

# La Era de "conversar" con los datos



CBINSIGHTS  
**Generative AI**  
**50**  
2023

**AI development tools**

<b>Foundation models &amp; APIs</b> AI21labs    ALEPH ALPHA cohere    contextual-ai Hugging Face    OpenAI	<b>Data curation</b> Cleanlab	<b>Model development &amp; fine-tuning</b> FIXIE    Lightning**    mindsdb
<b>Model observability</b> WHYLABS	<b>Vector database tech</b> LangChain    Weaviate    Zilliz	

**Cross-industry applications**

<b>Workplace knowledge management</b> glean Hebbia mem	<b>Synthetic voice</b> IIElevenLabs  <b>Design tools</b> Galileo AI Poly	<b>Search</b> Perplexity Twelve Labs vectara YOU	<b>Image &amp; video generation</b> MidJourney    synthesisia  <b>Sales &amp; customer support</b> AssemblyAI    tavus	<b>Code generation</b> diffblue replit warp	<b>AI assistants &amp; HMs</b> ADEPT ANTHROPIC Inflection
---	---	--	--	--	--

**Industry-specific**

<b>Education</b> Elio	<b>Gaming</b> convai	<b>Healthcare</b> AQEMIA    navina    SUBTLE MEDICAL	<b>Materials &amp; manufacturing</b> Cradle    TOROGRAM
<b>Legal</b> Harvey	<b>Construction</b> Augmenta	<b>Media &amp; entertainment</b> character.ai    descript    FLAMES    MOE    runway    wonder	

Note: Companies are private as of 7/25/23.

Imagen "Llama" extraída de ejemplos de en [midjourney.com/](https://midjourney.com/)

# Plataforma Inteligente

Nuestra apuesta



- Buscamos que los **prototipos de IA y sus resultados prometedores evolucionen**.
- El paso hacia la escala de **soluciones de IA industriales y confiables** es esencial para ofrecer valor a nuestros clientes y convertir la IA en una ventaja competitiva.
- Sopra Steria ha desarrollado la **Plataforma InnerData AI (basada en Red Hat OpenShift y el operador Open Data Hub)**, para **analistas y científicos de datos** capaces de producir y desplegar eficientemente modelos de IA.



- Facilitar el acceso y uso sencillo de un entorno completo orientado a datos e IA, desde la **ingestión de datos hasta la exposición compatible con la industria**.
- Servir como una base de software para una **arquitectura genérica de datahub automatizado** que permita la rápida realización de pruebas de concepto (POV) y servir como un núcleo básico para proyectos clave de clientes.
- Facilitar la creación progresiva de un **marketplace** de modelos de IA.



## Cookbooks

Recetas digitales para resolver problemas específicos de IA con un ejemplo empresarial



## Model Zoo

Modelos preentrenados con Datasets orientados a negocios



## Assets & Tools

Mostrar proyectos de IA realizados y permitir la capitalización de nuestros éxitos.



# Reference Architecture for AI on OpenShift sopra steria

## Artificial Intelligence & Machine Learning

<b>Model Training</b> <i>Tensorflow Training (TFJob)</i> <i>PyTorch Training</i> <i>Spark</i> <i>Katib</i>	<b>Model Serving</b> <i>Seldon Core</i> <i>Tensorflow Serving</i> <i>KFServing</i>	<b>Interactive Notebooks</b> <i>JupyterHub</i> <i>Notebook Server (Kubeflow)</i> <i>Elyra</i> <i>JupyterLab</i>	<b>ML Applications</b> <i>Open Data Hub AI Library</i>	<b>Business Intelligence</b> <i>Superset</i>
---	---	--	---	---

## Data Analysis

<b>Big Data Processing</b> <i>Spark</i> <i>Spark SQL</i> <i>Thrift</i> <i>Presto</i>	<b>Streaming</b> <i>Kafka Streams</i> <i>Elasticsearch</i>	<b>Data Exploration</b> <i>Hue</i> <i>Kibana</i>
--	---	---

## Metadata Management

*Hive Metastore*   *Metadata (Kubeflow)*

## Storage

<b>Data Lake</b> <i>Red Hat® Ceph Storage (Ceph)</i> <i>Red Hat® Openshift Container Storage (rook-ceph)</i>	<b>In-Memory</b> <i>Red Hat® Data Grid (Infinispan)</i>	<b>Relational Databases</b> <i>PostgreSQL</i> <i>MySQL</i>
--	--	---

## Data in Motion

*Red Hat® AMQ Streams (Kafka Strimzi)*   *Red Hat® Ceph S3 API*   *Red Hat® Fuse (Camel K)*   *Kafka Connect*   *Logstash*   *Fluentd*   *rsyslog*

## Security & Governance

*Red Hat® OpenShift OAuth*

*Red Hat® Single Sign-On (Keycloak)*

*Red Hat® Ceph Object Gateway*

*Red Hat® 3scale*

Data Steward

## Monitoring & Orchestration

*Prometheus*

*Grafana*

*Kubeflow Pipelines*

*Elyra*

*Kubeflow Experiments*

*Airflow*

*Argo Workflows*

*Openshift Pipelines (Tekton)*

*Argo CD*

DevOps Engineer

**Red Hat OpenShift**

**kubernetes**

**Red Hat Enterprise Linux**

**Red Hat Hybrid Cloud Management**



Versatile Data Storage



Dev environments



Pipelines & CI/CD



Perf. & Stat. Monitoring



Log monitoring & analysis



Reuse to build



Nuevo objetivo de la plataforma:

# Usar, entrenar, y evaluar LLMs

**1 Datos**  
Obtención de conjuntos de datos para entrenar LLMs, su contenido y problemas.



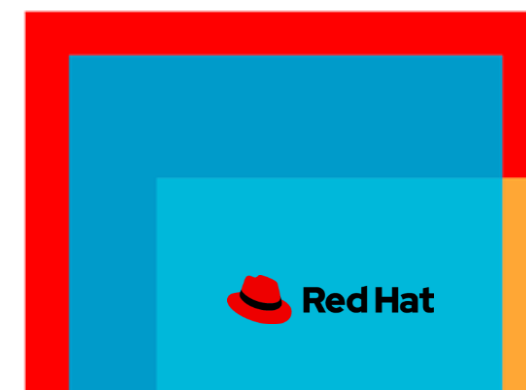
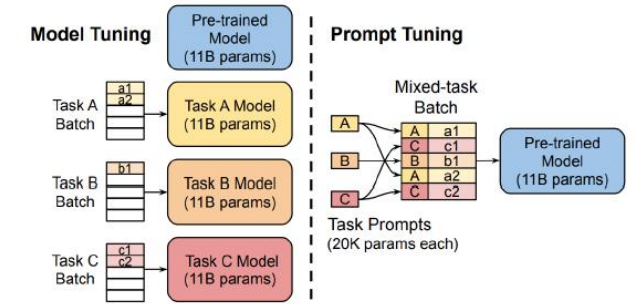
**2 Entrenamiento**  
Entrenamiento de una LLM desde 0, y fine-tuning para tareas específicas.



**3 Evaluación**  
Evaluación de LLMs y benchmarks utilizados en el estado del arte.

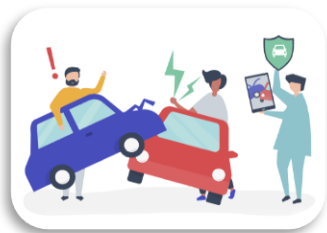


**4 Inferencia**  
Inferencia con LLMs, quantization, prompt-engineering, y casos de uso.



Nuevos retos en los casos de uso:

# Generar datos de entrenamiento



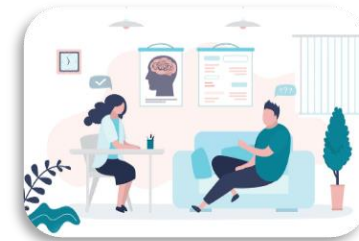
## Seguros

Crear perfiles de clientes simulados para predecir reclamaciones y calcular tarifas.



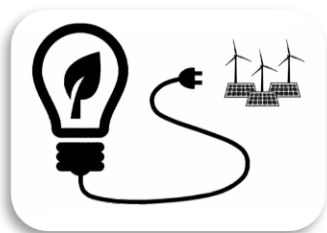
## Sociales y Demografía

Generar poblaciones virtuales y estudiar tendencias demográficas y la planificación urbana.



## Salud

Crear conjuntos de datos de pacientes simulados sin comprometer la privacidad de los pacientes.



## Energía y Medio Ambiente

Simulación de datos de sensores y redes eléctricas.



## Finanzas

Imitar el comportamiento de transacciones financieras para la detección de anomalías.

Otros...

Los LLMs no resuelven los viejos problemas:

# Problemas con los datos sintéticos

**Preocupaciones éticas y de seguridad:** uso indebido (generación de contenido engañoso o dañino).

**Sesgo:** reflejar y amplificar los sesgos presentes en sus datos de entrenamiento.

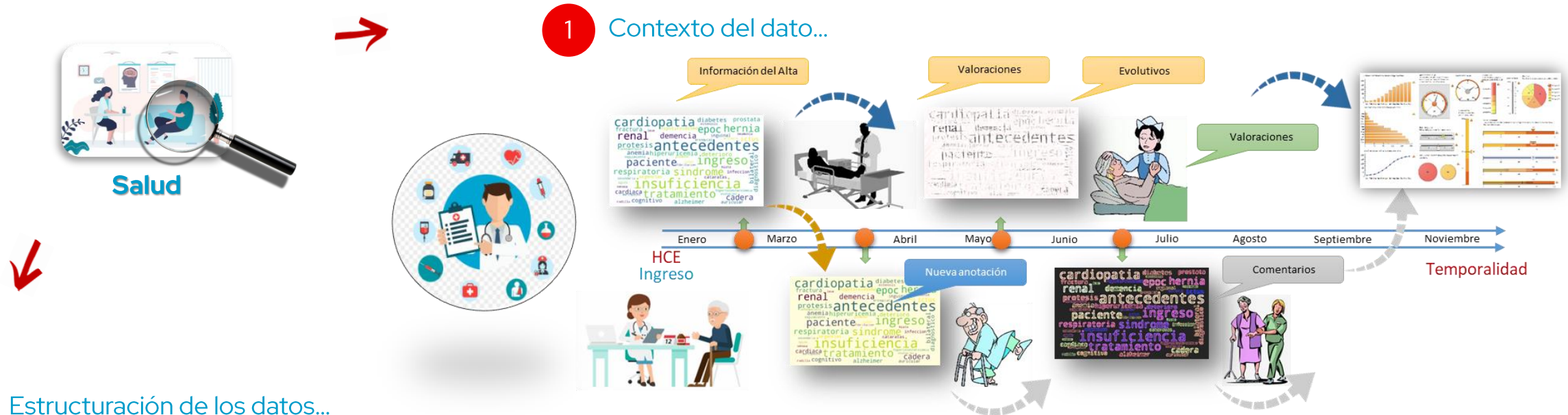
**Interpretabilidad y transparencia:** entre más grandes y complejos, más difícil es entender cómo toman sus decisiones, un problema para tareas donde la transparencia es importante.

**Generalización y solidez:** pueden tener dificultades con tareas que son ligeramente diferentes o con entradas que no han visto durante el entrenamiento.





# Problemas con el uso de datos sintéticos en contextos sensibles



## 2 Estructuración de los datos...



Imágenes autor: ponente

## Uso de técnicas de etiquetado débil sobre datos sintéticos

### Problema...

La necesidad de cantidades sustanciales de datos de entrenamiento etiquetados manualmente.

### Contexto...

La llegada del transfer learning ha aliviado enormemente este requisito, pero aún sigue siendo necesario generar más ejemplos etiquetados.

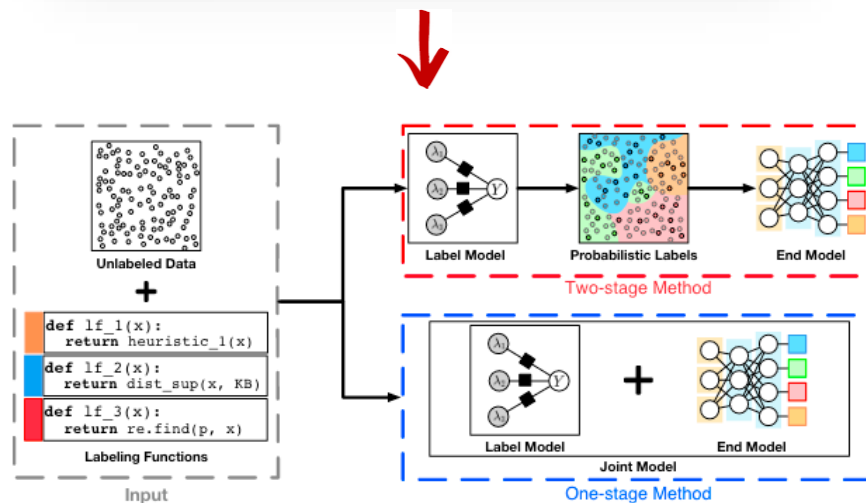


Imagen autor: <https://arxiv.org/abs/2109.11377>

### Weak Supervision

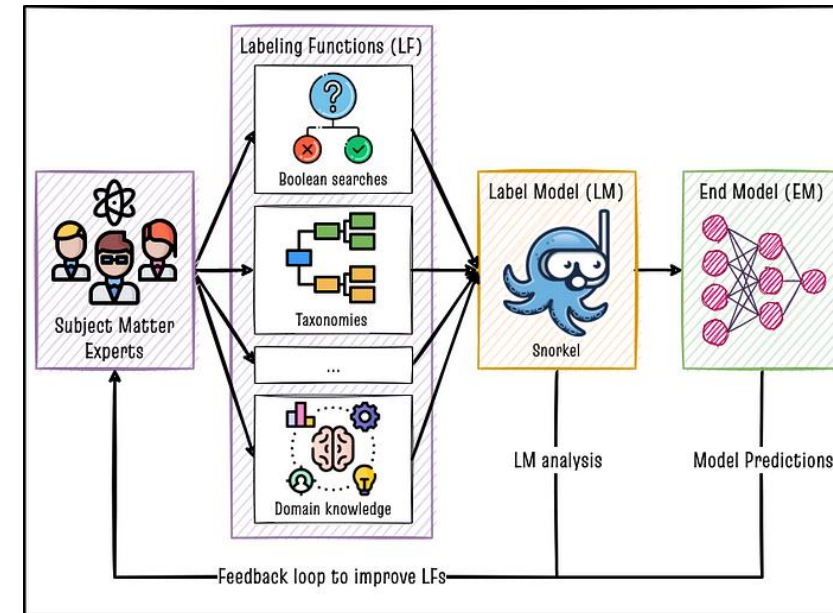


Imagen autor: [linkedin.com/in/marie-stephen-leo](https://www.linkedin.com/in/marie-stephen-leo)

## Formalización clásica de la generación de datos sintéticos

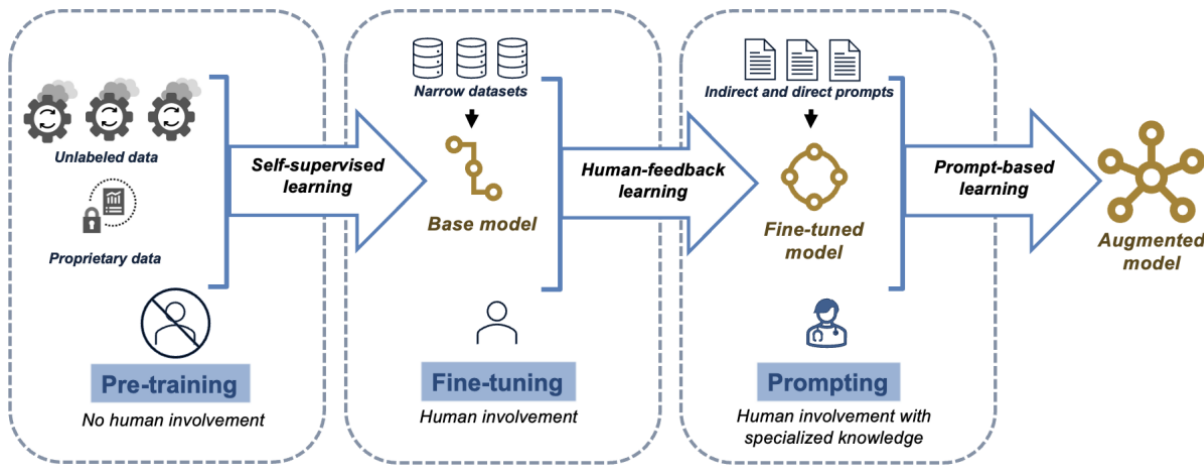


Imagen Autor: Jesutofunmi & Gui, Haiwen & Rezaei, Shawheen & Zou, James & Daneshjou, Roxana. (2023)

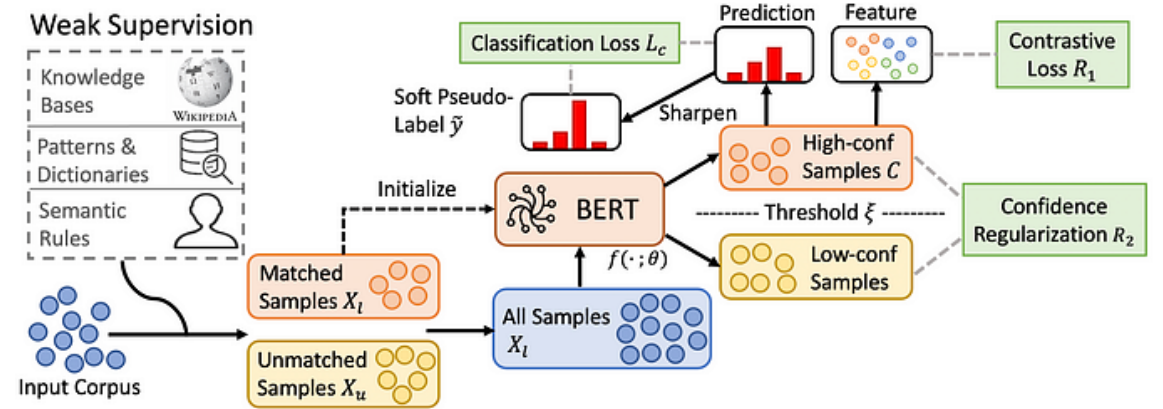
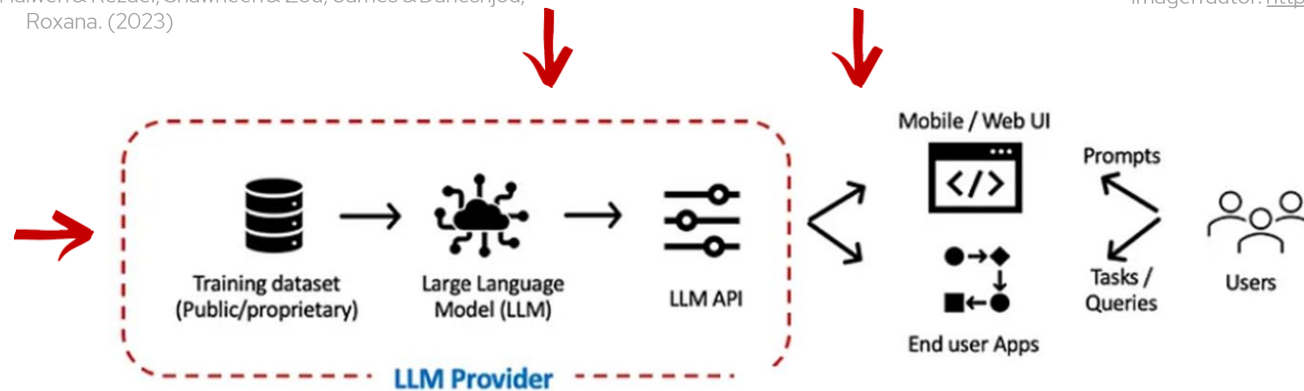


Imagen autor: <https://aclanthology.org/2021.naacl-main.84/>

Modelo Caja Negra



Tomado del artículo de [Debmalyas Biswar](#)

# LLMOps y el versionado de datos sintéticos

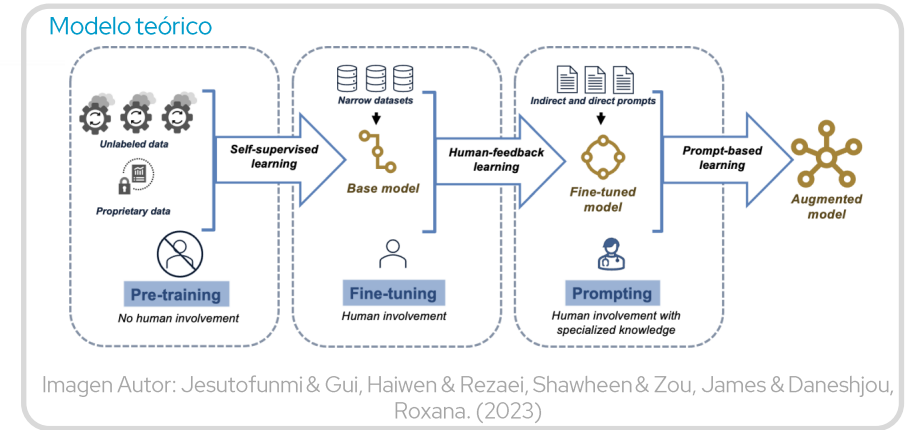
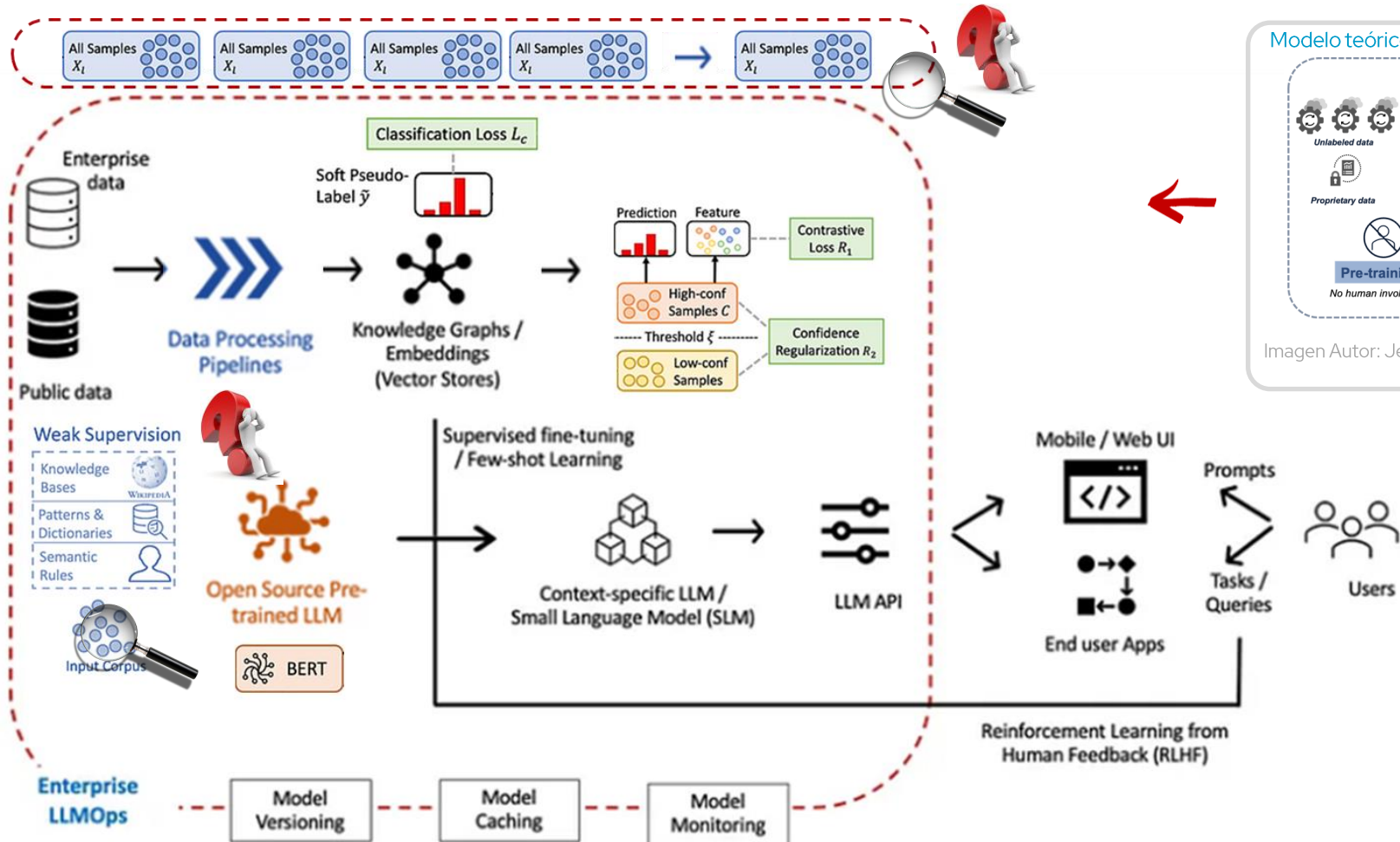


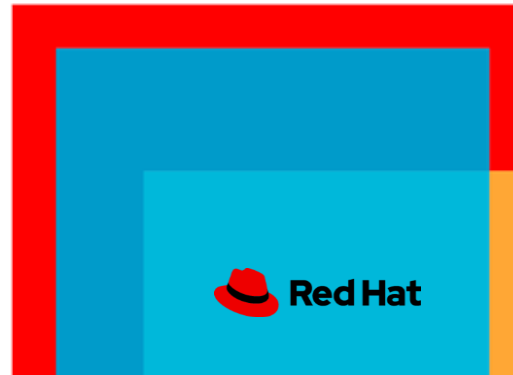
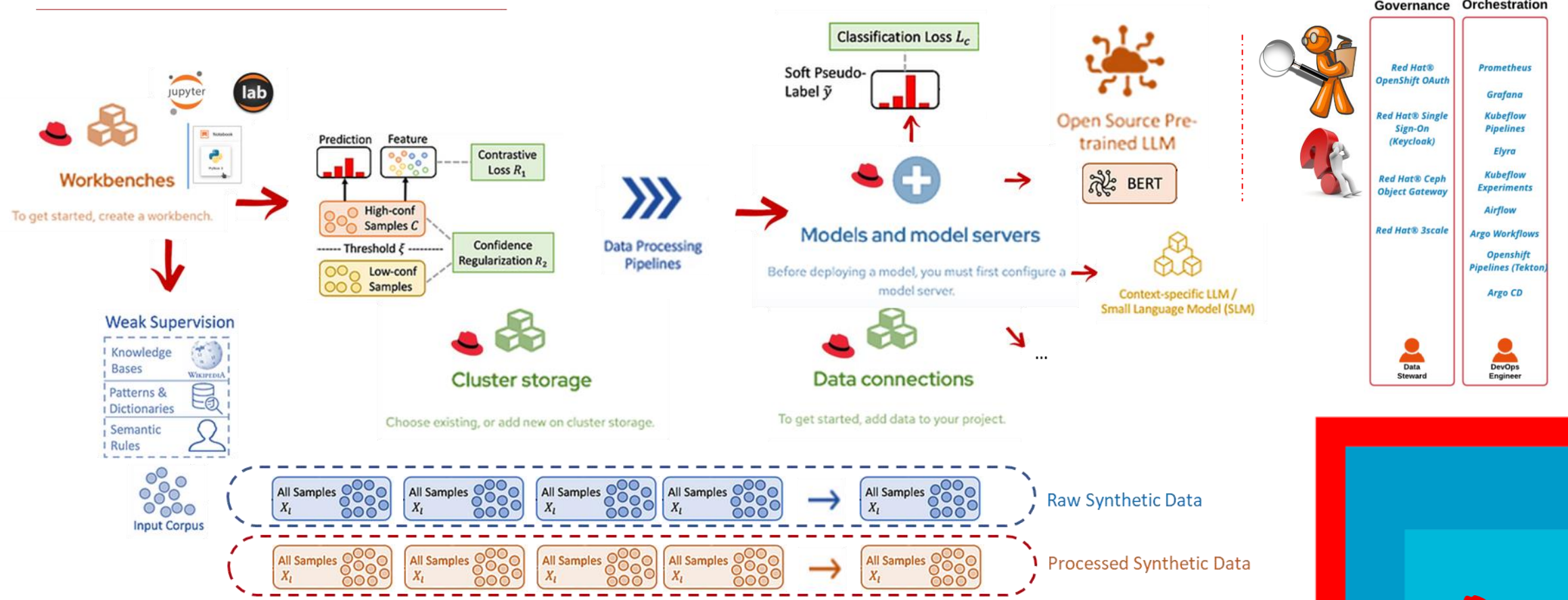
Imagen Autor: Jesutofunmi & Gui, Haiwen & Rezaei, Shawheen & Zou, James & Daneshjou, Roxana. (2023)

Inspirado en el artículo de [Debmalyas Biswar](#)

# Laboratorios para el versionado de datos sintéticos



**Red Hat**  
OpenShift Data Science



Red Hat  
**Summit**

**Connect**

Thank you



[linkedin.com/company/soprasteria](https://www.linkedin.com/company/soprasteria)



[facebook.com/SopraSteriaEspana/](https://www.facebook.com/SopraSteriaEspana/)



[youtube.com/@SopraSteriaGroup](https://www.youtube.com/@SopraSteriaGroup)



[twitter.com/SopraSteria\\_ES](https://twitter.com/SopraSteria_ES)