

Desafíos de la Inteligencia Artificial como servicio

Explicabilidad y Privacidad

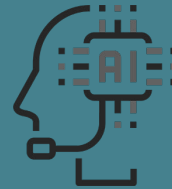
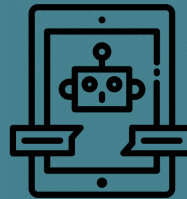
Eduardo Mangarelli
mangarelli@ort.edu.uy

Sergio Yovine
yovine@ort.edu.uy

- Grupo de investigación en IA
- Master en Big Data - <https://fi.ort.edu.uy/master-en-big-data>
 - Dip. Esp. en Analítica de Big Data (desde 2017)
 - Dip. Esp. en Inteligencia Artificial
- Más de 40 graduados (agosto 2019)
- Vinculación con empresas
- Cuerpo docente con académicos y profesionales
- Plataforma cloud de IA propia

AlaaS = (I + P + S)aaS

AI Services



**AI
Tools**



**Big Data
Storage**

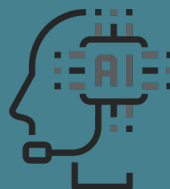
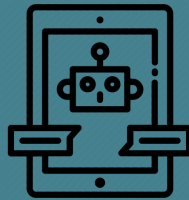


**Big Data
Processing**



Exposición de datos y modelos

AI Services



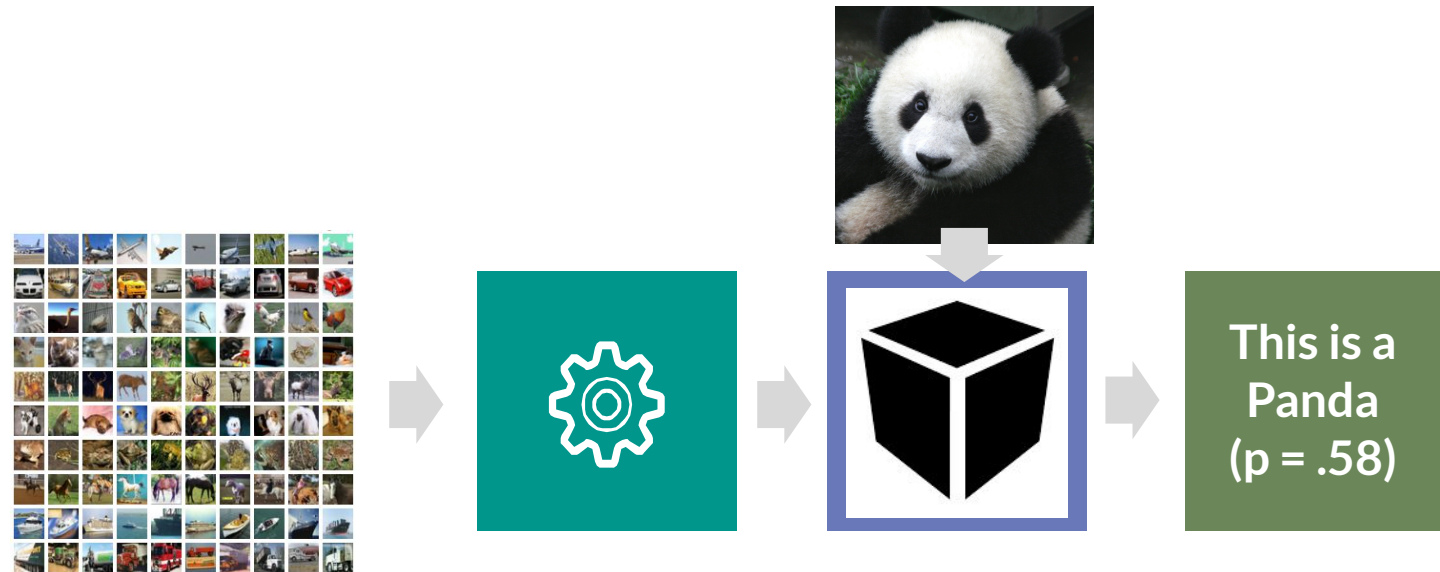
Big Data
Storage



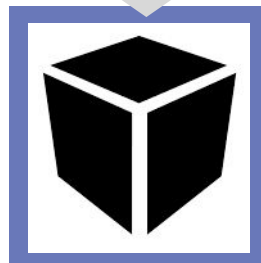
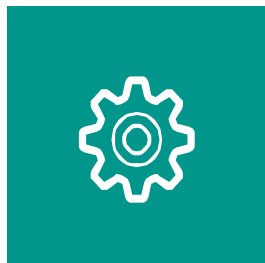
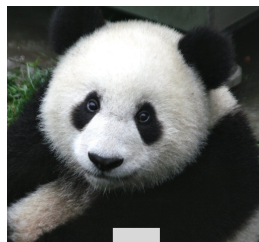
Big Data
Processing

AI
Tools

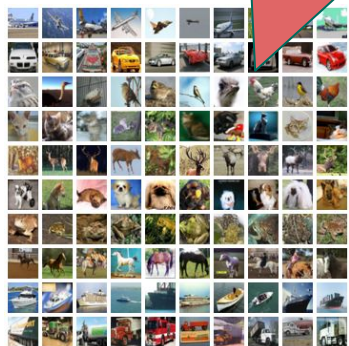




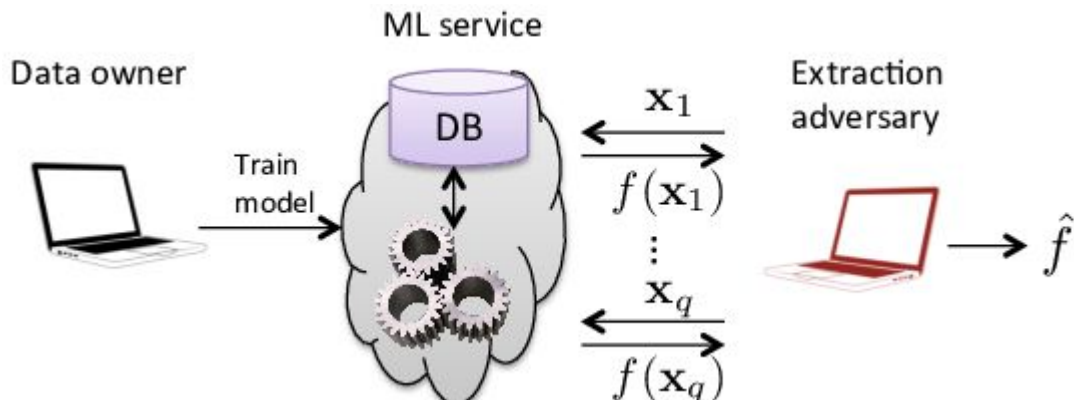
¿Mi foto
está
segura?!



This is a
Panda
($p = .58$)



Ataques a la privacidad



Extraer información sobre los datos de entrenamiento

Ataques a la privacidad

Dato de
entrenamiento

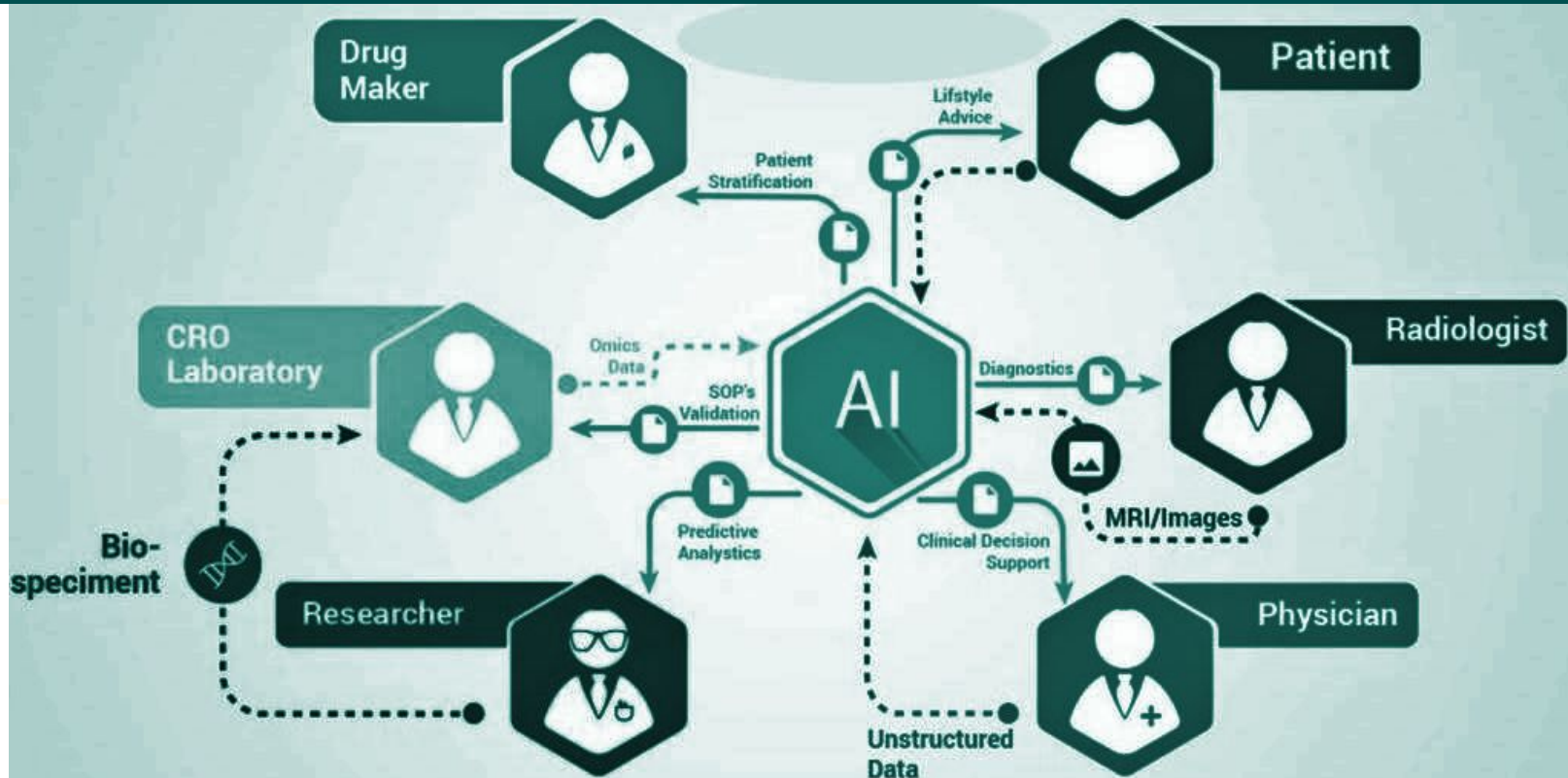


Dato
reconstruido



- Privacidad diferencial en el entrenamiento
 - Proteger los datos de entrenamiento
- Generación de datos sintéticos privatizados
 - Publicar datos con valor predictivo
 - Proteger los datos reales



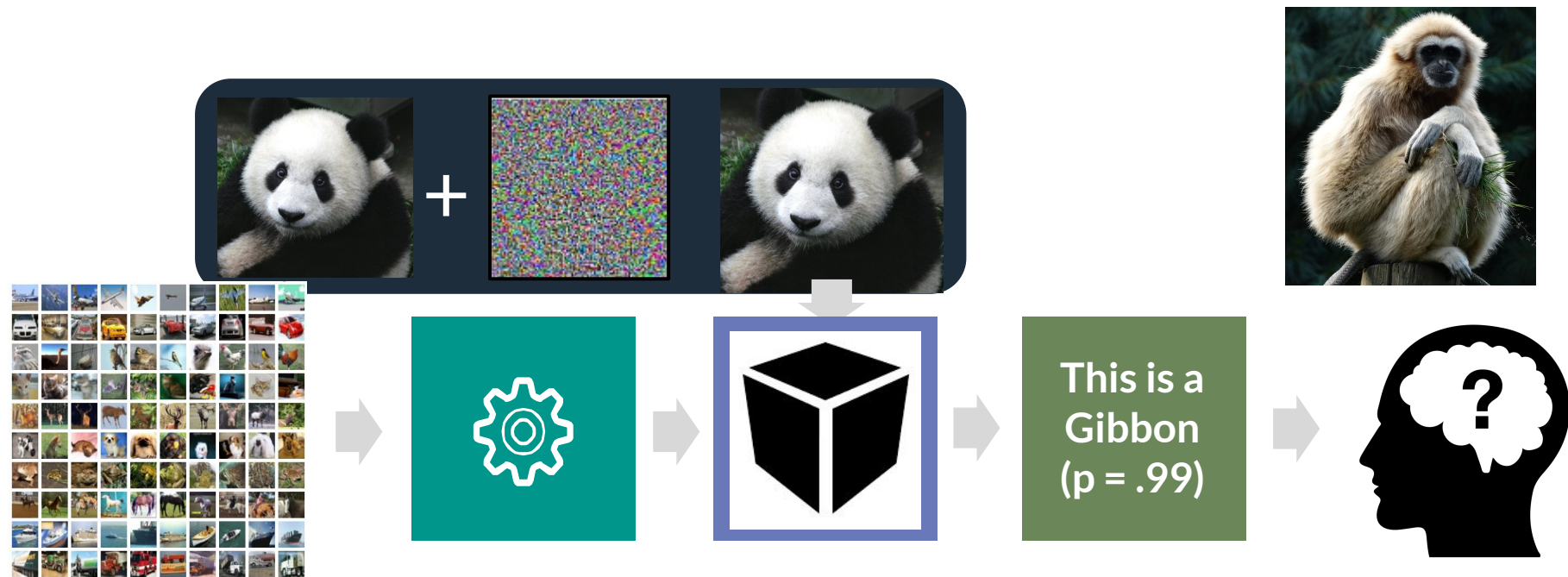


Exposición de datos y modelos



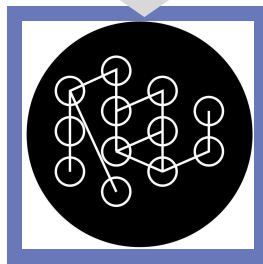
Interpretabilidad de los resultados





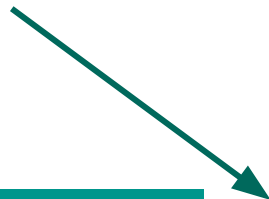
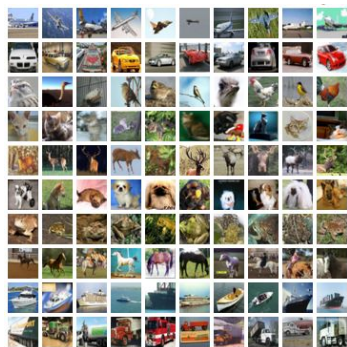
Modelo comprensible

Explicación



This is a
Cat

It has:
-fur
-whiskers
-claws



Soluciones para la explicabilidad

- Modelos interpretables
 - Bajo poder predictivo
- Explicar los modelos
 - Generar una derivación
 - Caja blanca
- Inducir un modelo interpretable
 - Aprendizaje activo (Q&A)
 - Caja negra



- Desafíos
 - Privacidad de los datos
 - Explicabilidad / Interpretabilidad / Verificación
- Proyectos en curso para abordarlos

