



# Plataformas distribuidas para análisis de grandes grafos

Dr. Cristian Mateos Díaz

CONICET / ISIS T A N-UNICEN

<http://www.exa.unicen.edu.ar/~cmateos>



**UNICEN**

UNIVERSIDAD NACIONAL DEL CENTRO  
DE LA PROVINCIA DE BUENOS AIRES

# Introducción

- Big data conduce a grandes volúmenes de datos, ¿pero cómo estructurarlos?
- Algunas respuestas desde la comunidad de Bases de Datos:
  - Key value databases (BigTable + GFS)
  - Wide-column databases (sharding + rows)
  - Document-oriented databases (XML)
  - *Graph databases*
- Dichos soportes se conocen como “**no SQL databases**”

# ¿Por qué grafos?

- Estructura versátil para modelado de datos
- Presencia en muchas aplicaciones
  - Redes sociales (Facebook, Twitter, Digg, Google+)
  - Scientometrics o bibliométricas
  - Arquitecturas de modelos biofísicos
  - ...
- Estas aplicaciones experimentan cambios de órdenes incesantes
  - Ejemplo: Más de 200 millones de usuarios activos, tweeting 400 millones de tweets por día; billones de relaciones seguidor/seguído

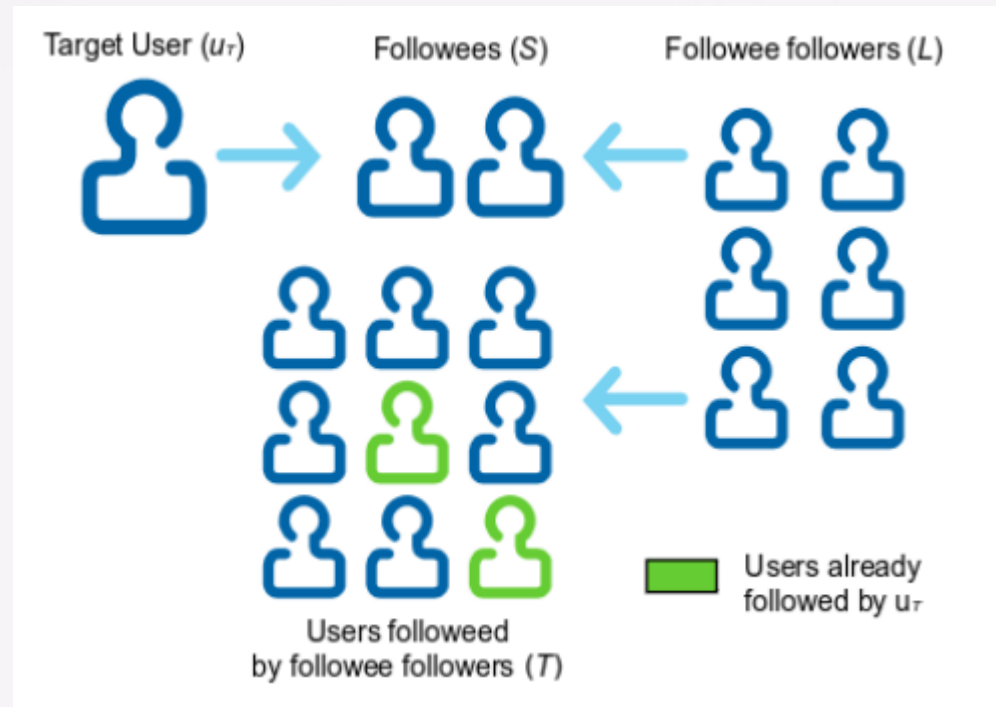
# Grafos y escalabilidad

- 30% de las bases de datos “no sql” son orientadas a grafos
  - Survey conducido durante 2013
  - Soporte para distribución y persistencia
  - Predominancia de bases clave-valor como alternativa
- Sin embargo, persisten problemas:
  - “Data sparseness”, que puede ser mitigado a nivel base de datos
  - Surgimiento de patrones de acceso por naturaleza ineficientes, que deben ser resueltos a nivel plataforma → no todas las particiones favorecen a todas las aplicaciones

# Grafos y escalabilidad: Twitter

- Un 77.9% de las relaciones son unidireccionales (celebridades)
- Sólo un 22.1% de las relaciones son recíprocas
  - Usado como fuente de información
- Un desafío que trae aparejado una oportunidad: **aplicaciones de recomendación**
  - WTF:[http://www.stanford.edu/~rezab/papers/wtf\\_overview.pdf](http://www.stanford.edu/~rezab/papers/wtf_overview.pdf))

# Grafos y escalabilidad: Twitter (continuación)



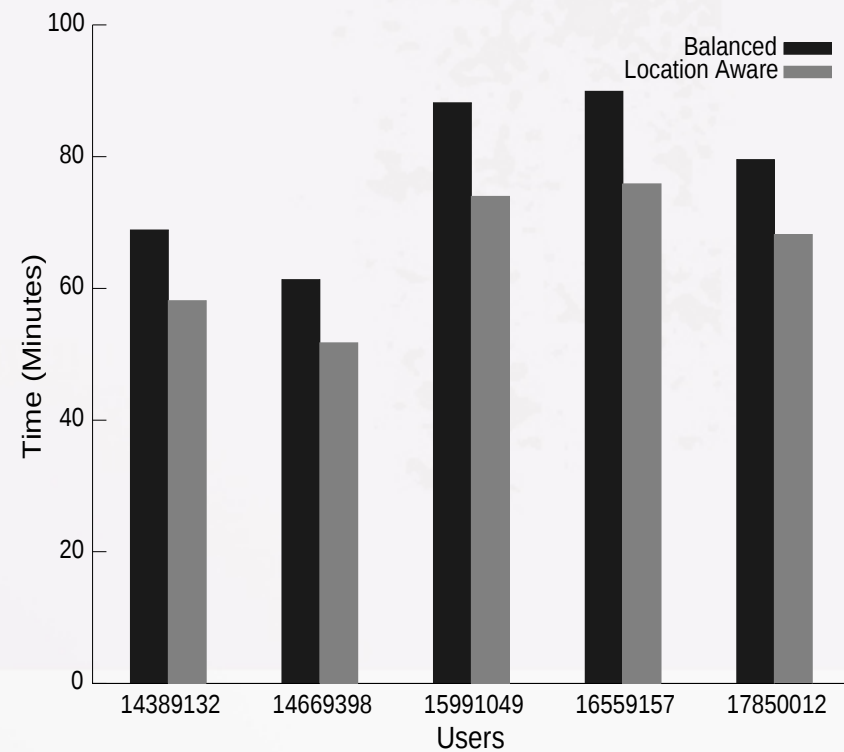
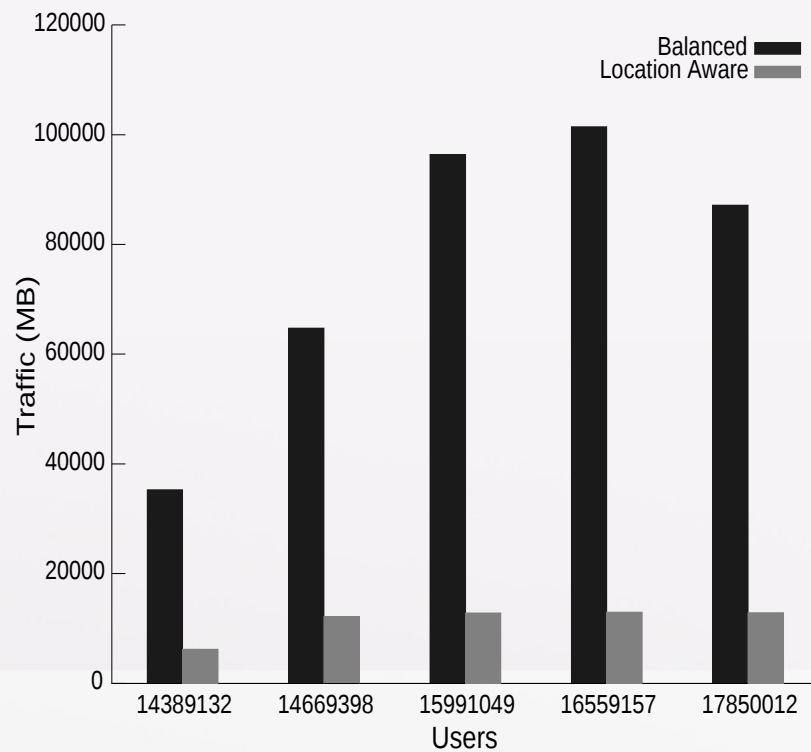
**Information Sources** → más seguidores que seguidos

**Friends** → usuarios con relaciones recíprocas

**Information Seekers** → usuarios con pocos tweets, y muchos seguidos

# Avances producidos

- Se están explorando heurísticas de **acceso a grafos**
- Experimentos llevados a cabo utilizando un dataset <http://an.kaist.ac.kr/traces/WWW2010.html> con 1.400 millones de relaciones
- Recomendaciones para 5 usuarios *information sources* ( $\text{followers}(u)/(\text{followers}(u)+\text{followees}(u)) \leq 0.5$ )



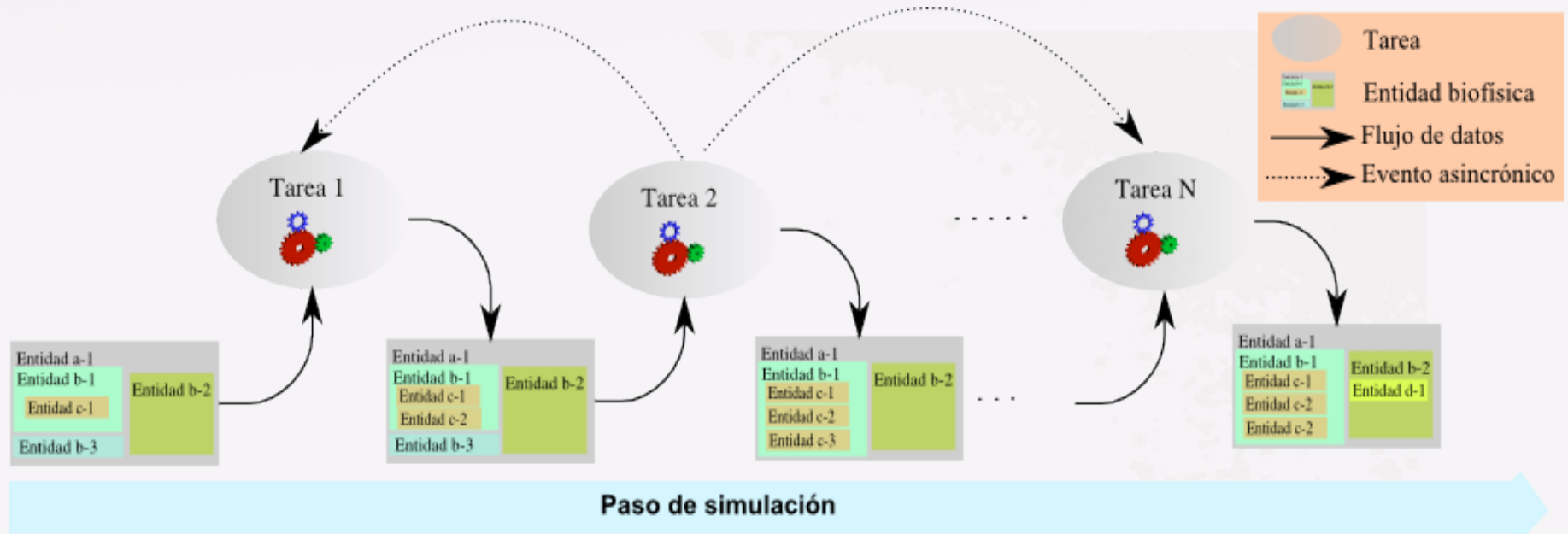


# Simulaciones agropecuarias

- Los sistemas de simulación agropecuaria son altamente dinámicos y complejos
  - Diversidad de entidades (ganado, pasto, ...)
  - Complejidad debido a la dimensión tiempo
  - Meteorología e incertidumbre económica
  - *Relaciones causales entre las entidades*
- ¡Sistemas computacionalmente demandantes!
  - Ejemplo: APSIM estudia la sustentabilidad productiva de cultivo (carbono, nitrógeno), 325 escenarios, 122 años → 30 años (secuencial), 10 días (paralelo)



# Simulaciones agropecuarias (continuación)



- Arquitecturas basadas en eventos (alta flexibilidad)
- Composición de entidades
- Paralelismo taskflow y dataflow
  - No tratados en *tandem* en la literatura...
  - Los datos son “ciudadanos de primera clase”

# Avances producidos

- Muchos sistemas de simulación agropecuaria adhieren a esta **estructura** común
- Predominancia de soluciones ad-hoc para escalabilidad → necesidad de una solución más general
- Se está estudiando la “gridificación” de estas aplicaciones
  - Desafíos: Sparseness propios del dominio y patrones de cómputo híbridos (taskflow + dataflow)
  - Objeto de estudio: Simugan

# Feature selection

- Problema actual en la comunidad de Machine Learning
  - Usada en tagging automático (ejemplo: Digg y similares)
- Eficiencia implica procesar matrices de millones x millones mediante nuevos algoritmos distribuidos
  - Problemas de sparseness y patrones de acceso
- Trabajando en primeros prototipos publicables (bibliotecas de álgebra lineal, GPUs no parecen ser suficientes!)

# Conclusiones

- Nuestra visión: Son necesarias formas de acceso eficiente a los datos
  - Enfoque común; MapReduce + BigTable
  - ¿Qué proveer? Modelos de programación, planificadores, cachés, y más...
- Experiencia académica en generación de tecnologías para sistemas distribuidos
  - PAE-PICT-2311 (Web Services)
  - PAE-PICT-2312 (Grids computacionales)
  - PICT-2012-0045 (Grids móviles)
  - PIP/PICT solicitados en la temática

**¡Gracias!**

