# GraphX

Amber, 2017.10.13

# Agenda

1.Spark + Scala + Jupyter
2.Spark GraphX
3.Hands On

# Spark + Scala + Jupyter

```
$ git clone https://github.com/jupyter-scala/jupyter-scala.git
$ cd jupyter-scala/
$ sh jupyter-scala


Run jupyter console with this kernel with
   jupyter console --kernel scala
Use this kernel from Jupyter notebook, running
   jupyter notebook
and selecting the "Scala" kernel.


$ jupyter kernelspec list

Available kernels:
  python2   /Users/…/anaconda2/lib/python2.7/site-packages/ipykernel/resources
  scala     /Users/…/Library/Jupyter/kernels/scala
```
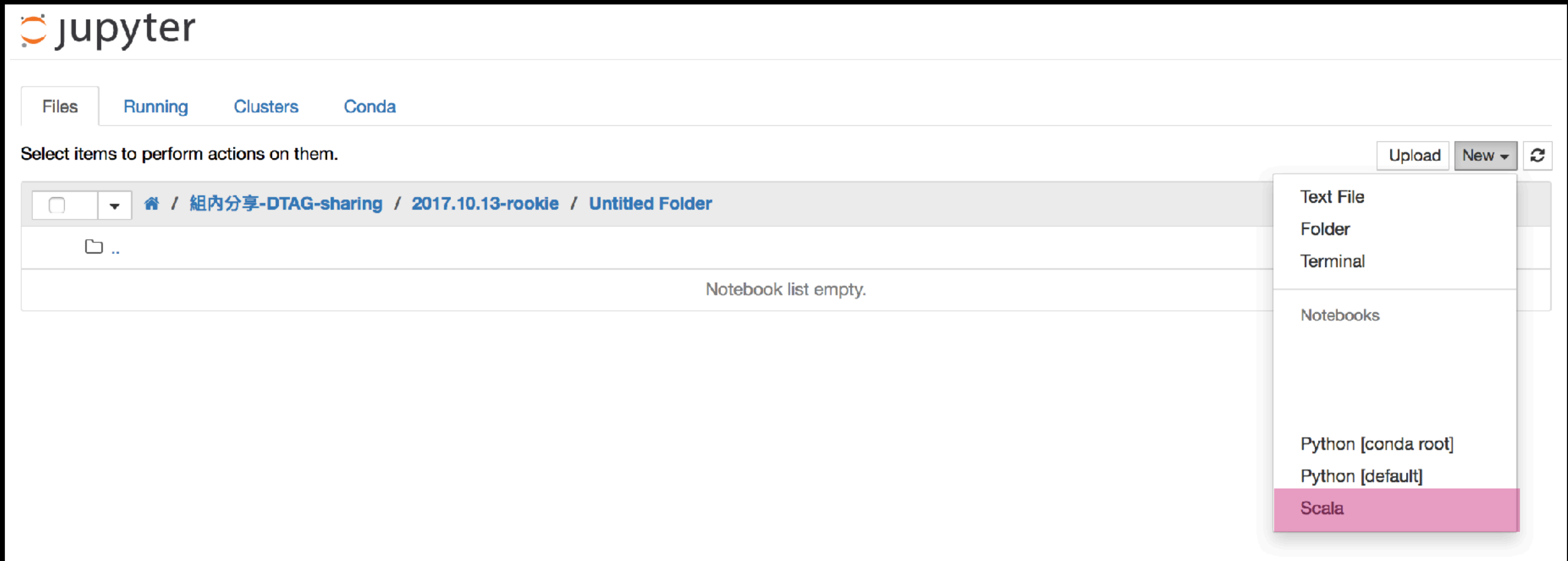
# Spark + Scala + Jupyter

`$  jupyter notebook`



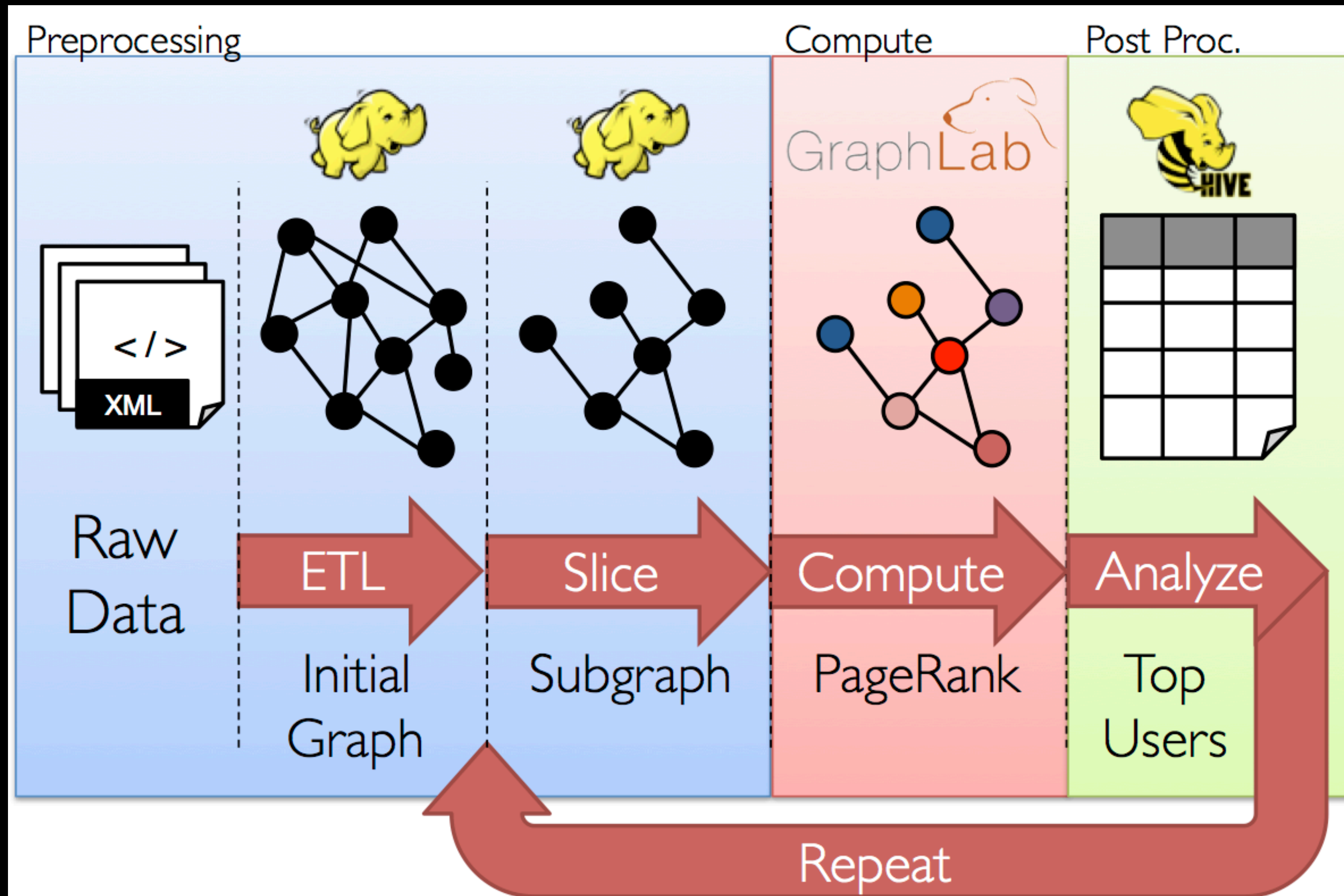線上管理頁面：**http://localhost:4040/jobs/**

# Spark GraphX

# Spark GraphX 介紹

GraphX 透過引入 Resilient Distributed Property Graph：一種帶有頂點和邊屬性的有向多重圖，來擴展Spark RDD。

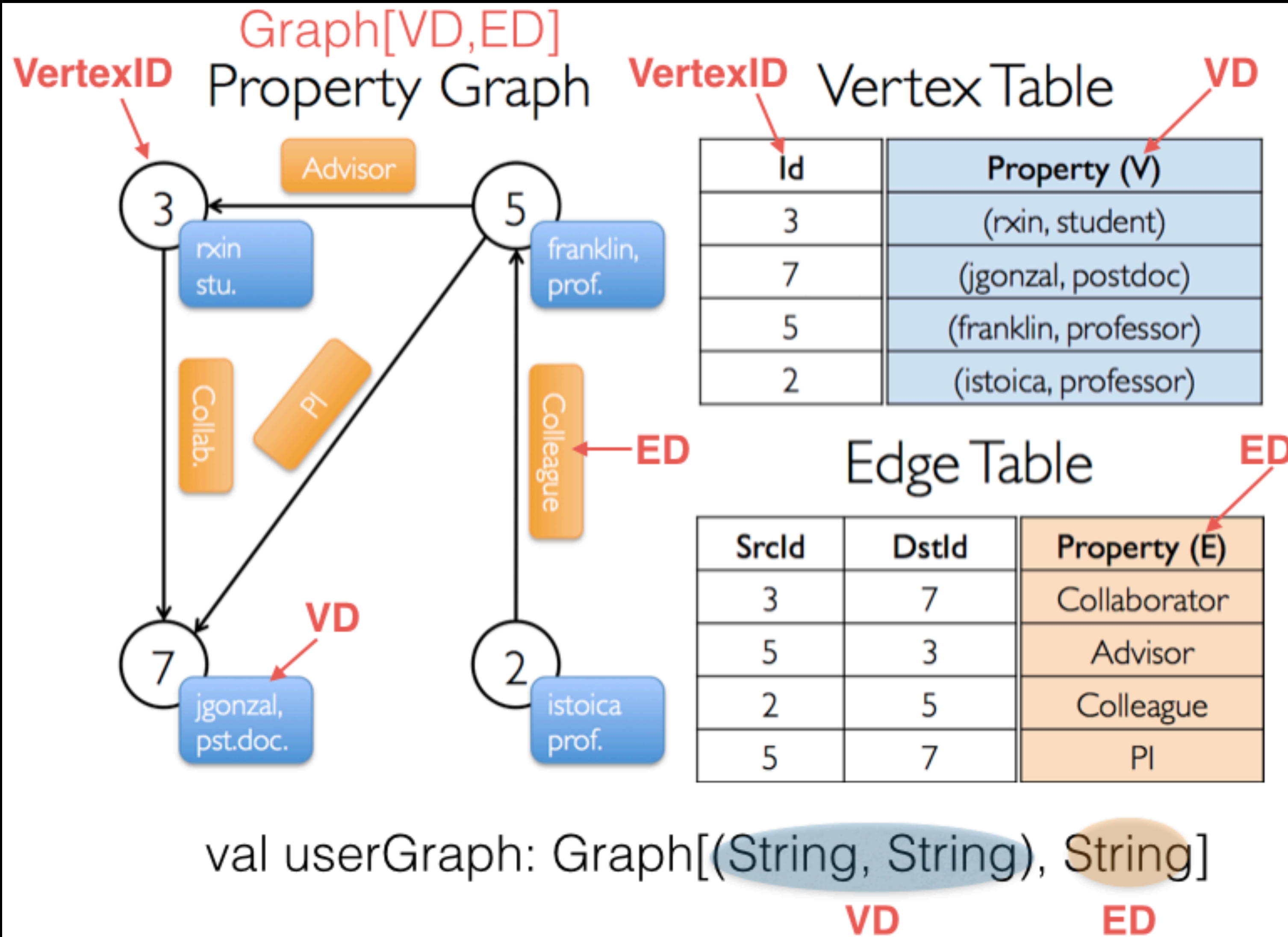為了支援圖形的運算，GraphX 擁有基本運算子和 Pregel API 的優化。
- 例如：subGraph、joinVertices、aggregateMessages

# Spark GraphX 介紹



**Graphx允許使用者將資料視為一個圖形和集合（例如：RDDs），而不需要任何的資料搬移和複製。**

# **Spark GraphX 資料結構**



Graph[VD,ED]

Property Graph

VertexID
VD

**Vertex Table**

| Id | Property (V) |
|----|--------------|
| 3 | (rxin, student) |
| 7 | (jgonzal, postdoc) |
| 5 | (franklin, professor) |
| 2 | (istoica, professor) |

**Edge Table**

| SrcId | DstId | Property (E) |
|-------|-------|--------------|
| 3 | 7 | Collaborator |
| 5 | 3 | Advisor |
| 2 | 5 | Colleague |
| 5 | 7 | PI |

ED
VD

val userGraph: Graph[(String, String), String]

VD        ED

**VertexRDD[VD] = RDD[(VertexId,VD)]**
VertexRDD[(String,String)] = RDD[(VertexId, (String,String))]

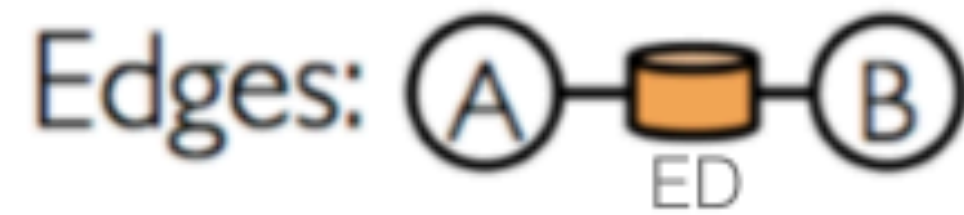**EdgeRDD[ED] = RDD[Edge[ED]]**
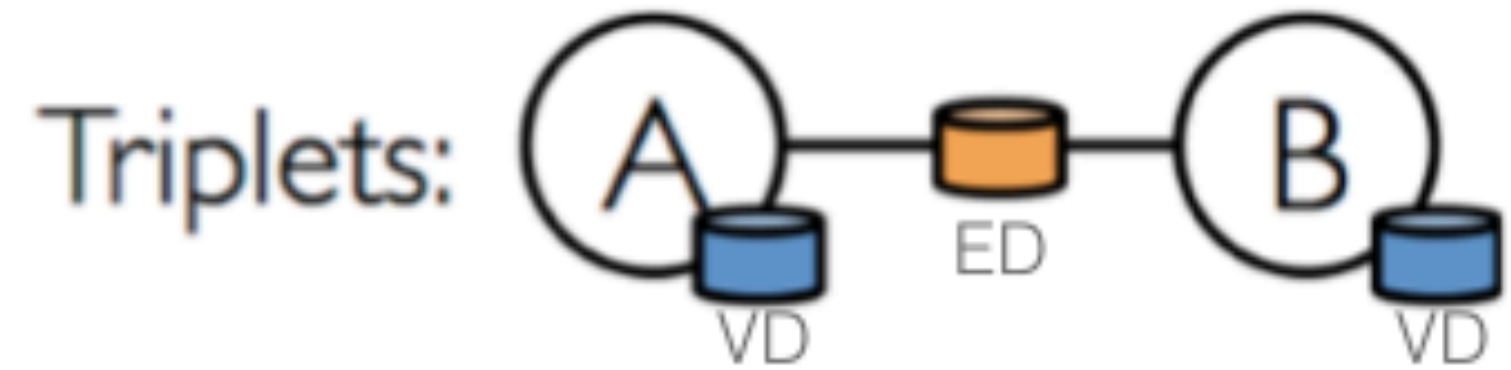EdgeRDD[String] = RDD[Edge[String]]

# Spark GraphX 資料結構



頂點視圖

邊視圖

三元組視圖

Vertices: A VD / B VD

RDD[(VertexId,VD)]

RDD[(VertexId, (String,String))]

+

Edges: A — ED — B

RDD[Edge[ED]]

RDD[Edge[String]]

=

Triplets: A — ED — B / VD ED VD

RDD[EdgeTriplet[VD, ED]]
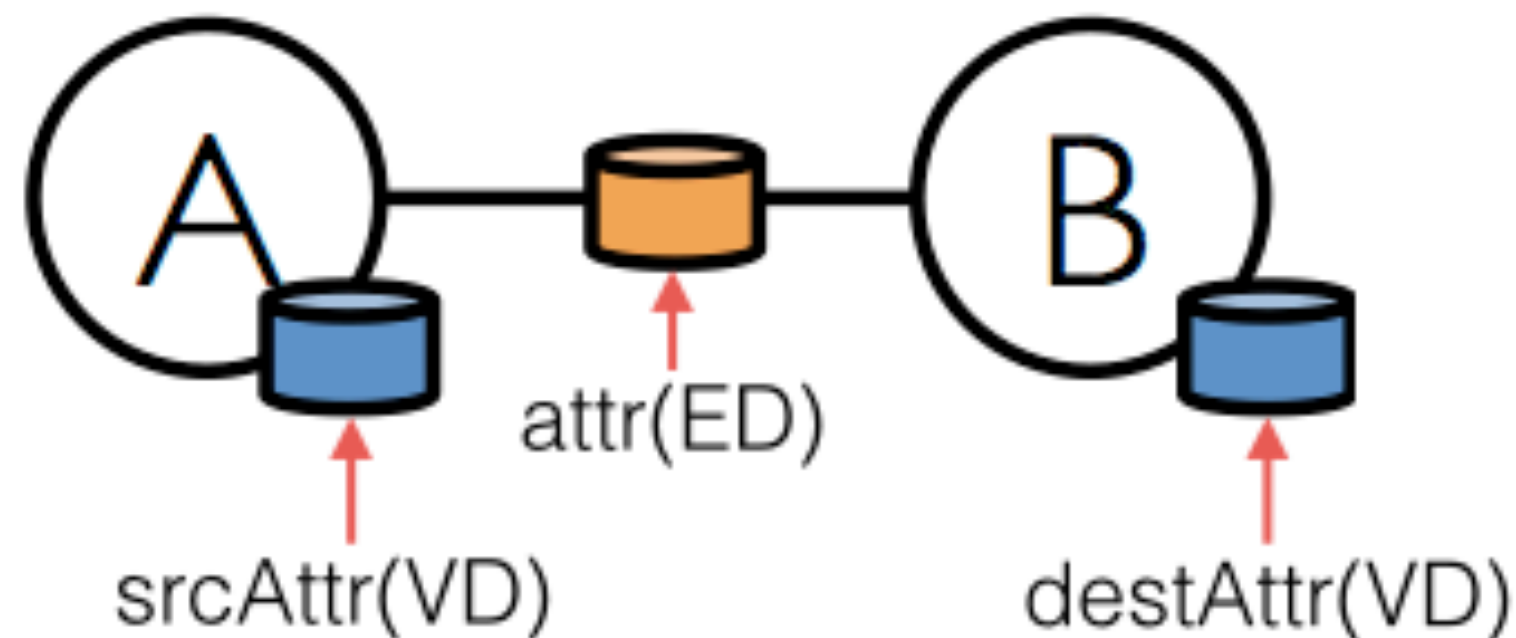
RDD[EdgeTriplet[(String,String), String]]

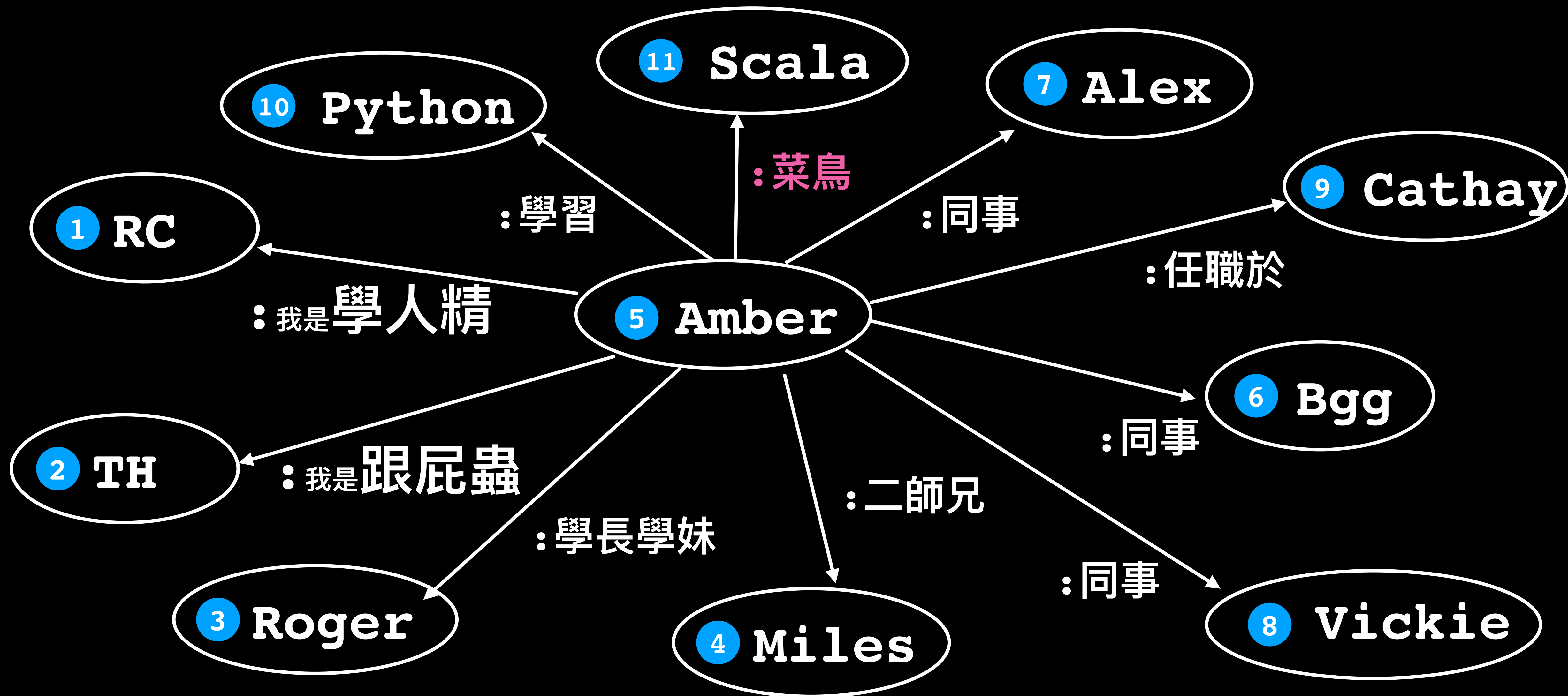| SQL | src:[id,attr]<br>dest:[id,attr] | e:[attr,srcId,destId] | SELECT src.id, dst.id, src.attr, e.attr, dst.attr<br>FROM edges AS e<br>LEFT JOIN vertices AS src, vertices AS dst<br>ON e.srcId = src.Id AND e.dstId = dst.Id |

Triplets: A — attr(ED) — B / srcAttr(VD) destAttr(VD)

EdgeTriplet[(String,String), String]

| attribute | type | example |
|-----------|------|---------|
| srcAttr | VD | (String,String) |
| destAttr | VD | (String,String) |
| attr | ED | String |

# 先來一點參與感

# About Amber

下集待續

# Agenda

1. Pregel API
2. Graph Algorithm
   - PageRank
   - Connected Components
3. Case Study