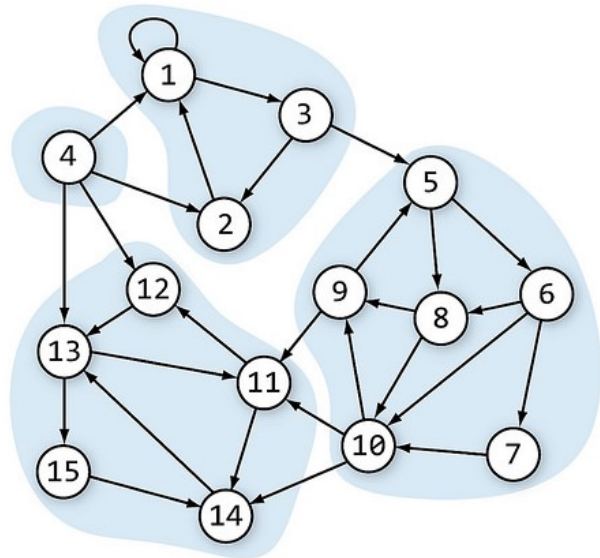


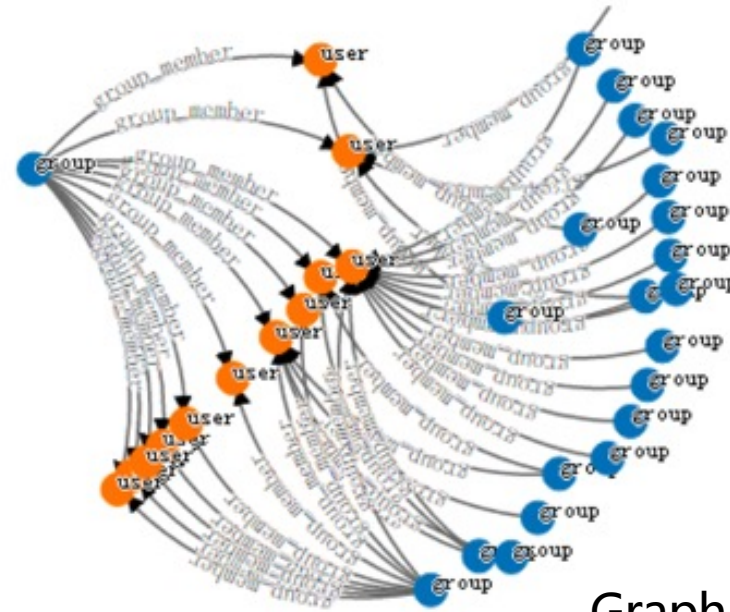
Apache GraphAr: An Open-Source Standard File Format for Graph Data Storage and Retrieval

Jingbo Xu
Institute for Intelligent Computing
Alibaba

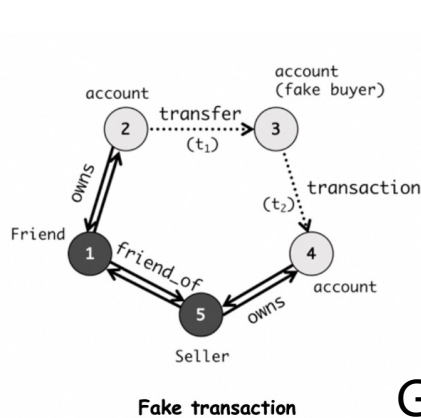
Diverse Graph Computing Workloads



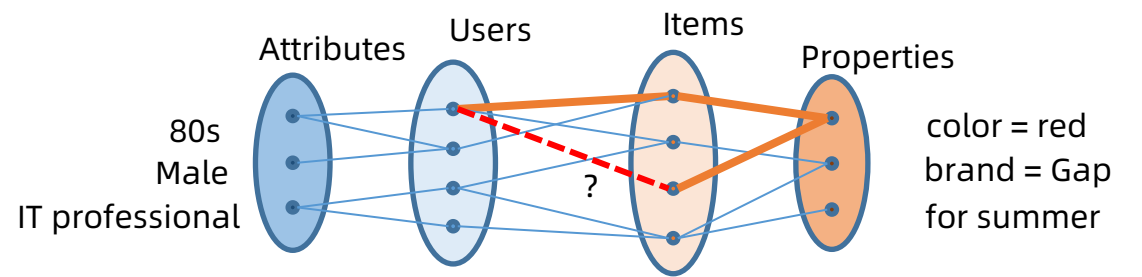
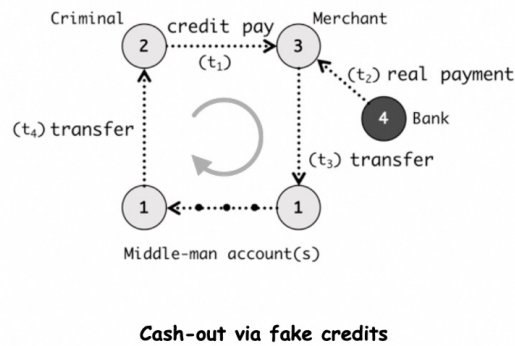
Iterative graph analytics, e.g., PageRank, LPA...



Graph querying/traversal

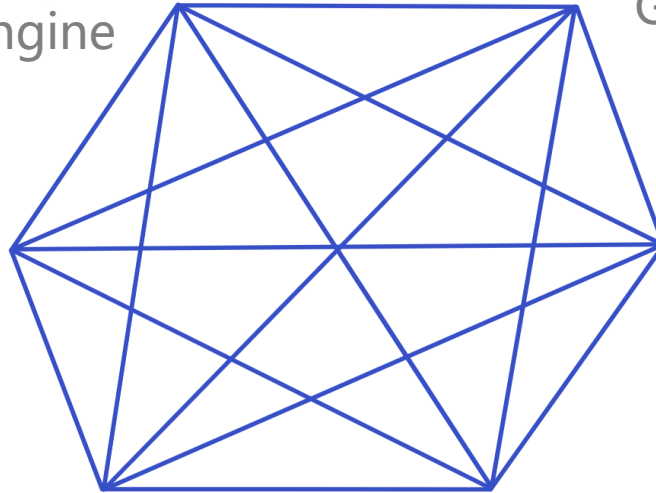
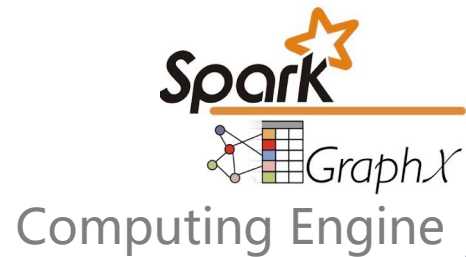


Graph pattern matching

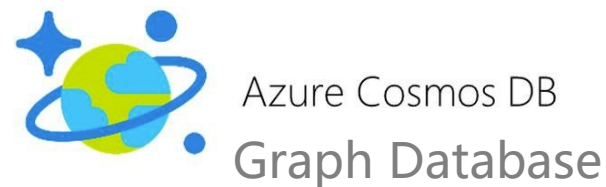


Various Graph Systems

- Graph DBs, GNN Frameworks, Graph Engines, KBs....
- **How do they share graph data?**



Amazon Neptune



What is GraphAr



GraphAr

- Short for **Graph Archive**
- An open source, standard data file format for graph data storage & retrieval.
- Designed for
 - Serving as a standardized file format for importing, exporting and archiving of the graph data which can be used by diverse existing systems, reducing the overhead when various systems co-work.
 - As a direct data source for (out-of-core) graph processing applications.

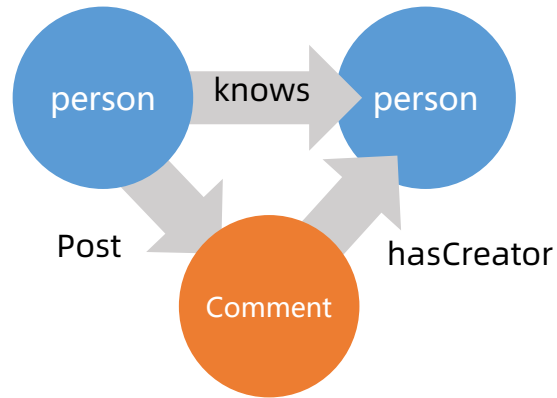
Key Features



- **Apache incubating**, open-source and vendor neutral
- The file format supports the **property graphs** and different representations for the graph topology (COO, CSR and CSC).
- It is compatible with existing widely-used file formats including **ORC, Parquet** (and less ideally CSV).
- **Apache Spark** can be utilized to generate, load and transform GraphAr files.
- It is convenient for use in a variety of **single-machine/distributed** graph processing systems, databases, and other downstream computing tasks.
- It enables users to conveniently perform **operations without modifying the payload files**, such as appending new vertices, adding new properties, or constructing a new graph with a set of selected vertices and edges.

Design of GraphAr: A Running Example

- A property graph and its vertices



internal vid	id	firstName	lastName	gender
0	933	Mahinda	Perera	male
1	6597069767117	Eli	Peretz	female
2	10995116278700	Joseph	Anderson	female
...
903	32985348834100	Bruno	Oliveira	male

Logical Table of Vertices

internal vid	id
0	933
1	6597069767117
...	...
499	13194139534267

./vertex/person/id/chunk0

internal vid	firstName	lastName	gender
0	Mahinda	Perera	male
1	Eli	Peretz	female
...
499	Asha-Rose	Chung	male

./vertex/person/firstName_lastName_gender/chunk0

internal vid	id
500	15393162788965
501	15393162789614
...	...
903	32985348834100

./vertex/person/id/chunk1

internal vid	firstName	lastName	gender
500	Hans	Becker	male
501	Adi	Cohen	female
...
903	Bruno	Oliveira	male

./vertex/person/firstName_lastName_gender/chunk1

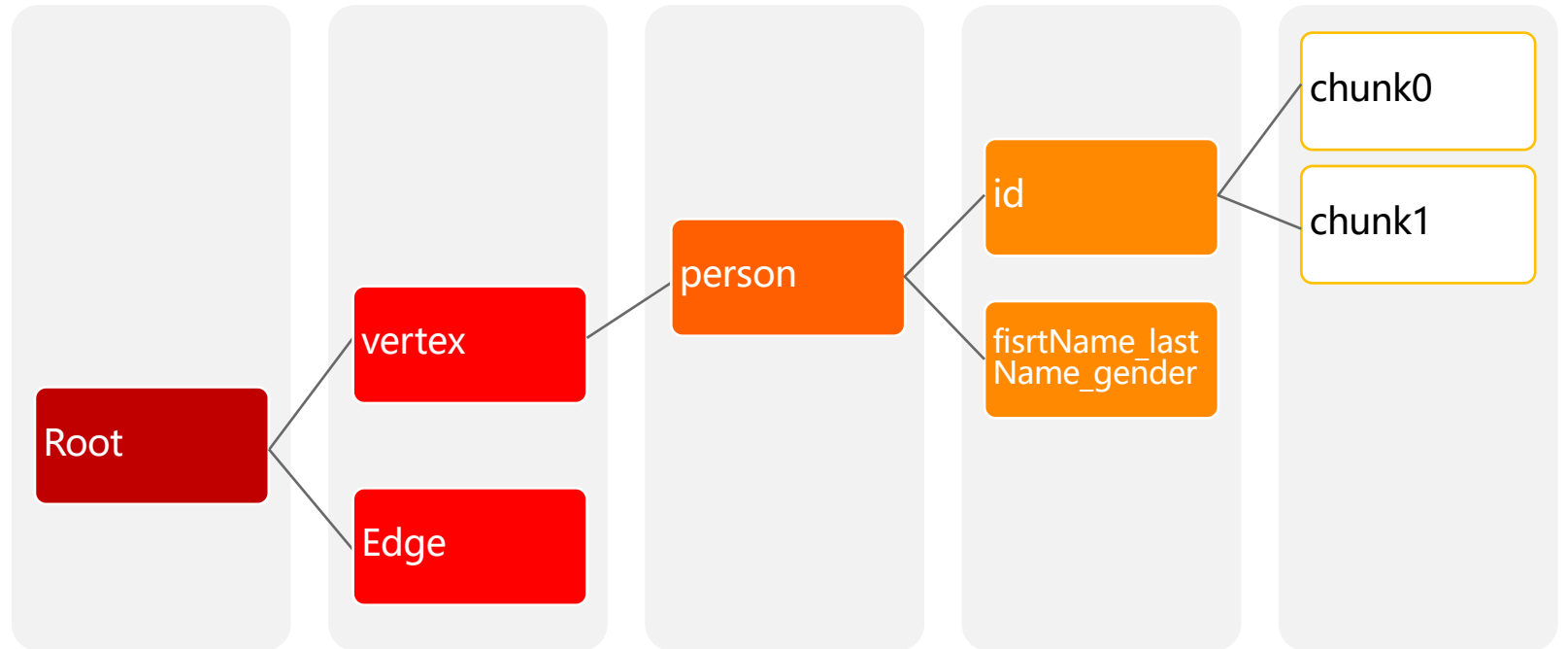
Physical Table of Vertices
Chunk size = 500

Vertices in GraphAr: Metadata and Payload Files

Metadata person.vertex.yml

```
1 label: person
2 chunk_size: 100
3 prefix: vertex/person/
4 property_groups:
5   - properties:
6     - name: id
7       data_type: int64
8       is_primary: true
9     prefix: id/
10    file_type: csv
11   - properties:
12     - name: firstName
13       data_type: string
14       is_primary: false
15     - name: lastName
16       data_type: string
17       is_primary: false
18     - name: gender
19       data_type: string
20       is_primary: false
21     prefix: firstName.lastName_gender/
22     file_type: csv
23 version: gar/v1
```

Payload files



Edges in GraphAr

Physical Table of Edges

Logical Table of Edges

internal vid	offset	source	destination	creationDate
0	0	0	87	2010-07-30T15:19:53.298+0000
1	3	0	849	2010-10-27T02:33:06.288+0000
...
901	6625	1	538	2012-04-12T13:56:58.931+0000
902	6626	1	539	2011-03-25T14:49:23.134+0000
903	6626
	6626	1	58	2011-06-10T17:47:19.432+0000
	6626
	6626	900	826	2012-08-21T00:02:08.530+0000
	6626	901	252	2012-08-13T10:11:20.606+0000

edges of vertex chunk 0

internal vid	offset
0	0
1	3
...	...
499	3772
	3778

source	destination
0	87
...	...
164	829

`./edge/person_knows_person/ordered_by_source/adj_list/part0/chunk0`

source	destination
164	30
...	...
269	565

`./edge/person_knows_person/ordered_by_source/adj_list/part0/chunk1`

source	destination
269	321
...	...
499	628

`./edge/person_knows_person/ordered_by_source/adj_list/part0/chunk2`

edges of vertex chunk 1

internal vid	offset
500	0
501	1
...	...
903	2848
	2848

source	destination
500	623
...	...
637	704

`./edge/person_knows_person/ordered_by_source/adj_list/part1/chunk0`

source	destination
638	375
...	...
793	884

`./edge/person_knows_person/ordered_by_source/adj_list/part1/chunk1`

source	destination
793	662
...	...
901	252

`./edge/person_knows_person/ordered_by_source/adj_list/part1/chunk2`

creationDate
2012-04-21T19:08:41.647+0000
...
2012-08-10T02:49:19.288+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part0/chunk0`

creationDate
2012-06-26T02:41:08.148+0000
...
2012-03-10T06:07:41.141+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part0/chunk1`

creationDate
2012-02-19T06:42:02.399+0000
...
2012-08-13T10:11:20.606+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part0/chunk2`

creationDate
2010-07-30T15:19:53.298+0000
...
2010-06-11T19:23:42.146+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part1/chunk0`

creationDate
2010-05-16T17:41:47.623+0000
...
2011-12-22T17:56:13.491+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part1/chunk1`

creationDate
2012-01-04T13:29:11.784+0000
...
2012-08-03T01:00:51.312+0000

`./edge/person_knows_person/ordered_by_source/creationDate/part1/chunk2`

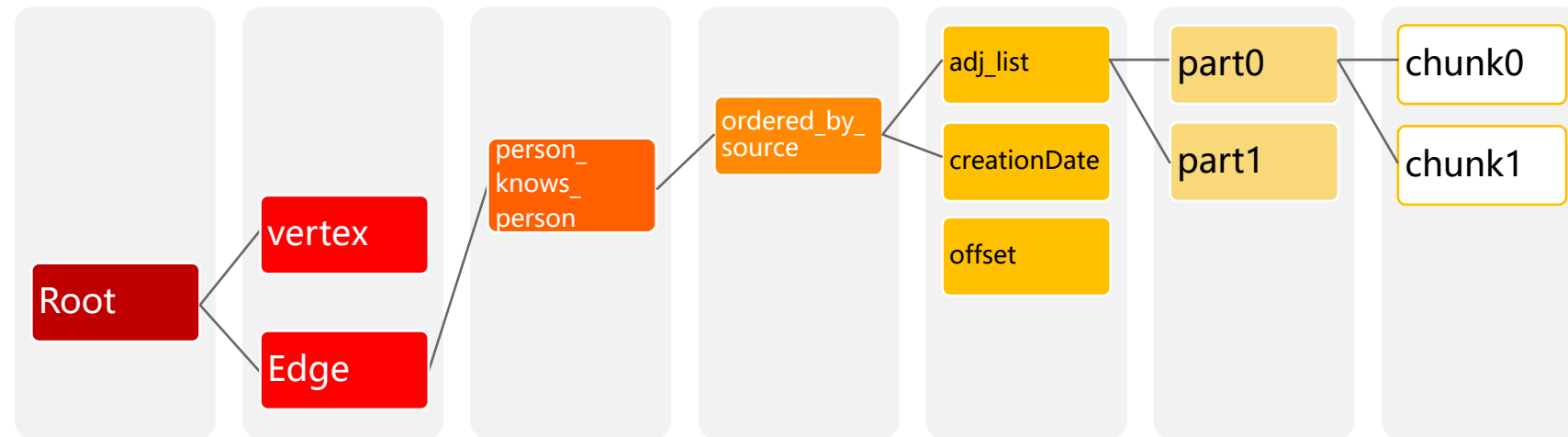
Edges in GraphAr: Metadata and Payload Files

Metadata

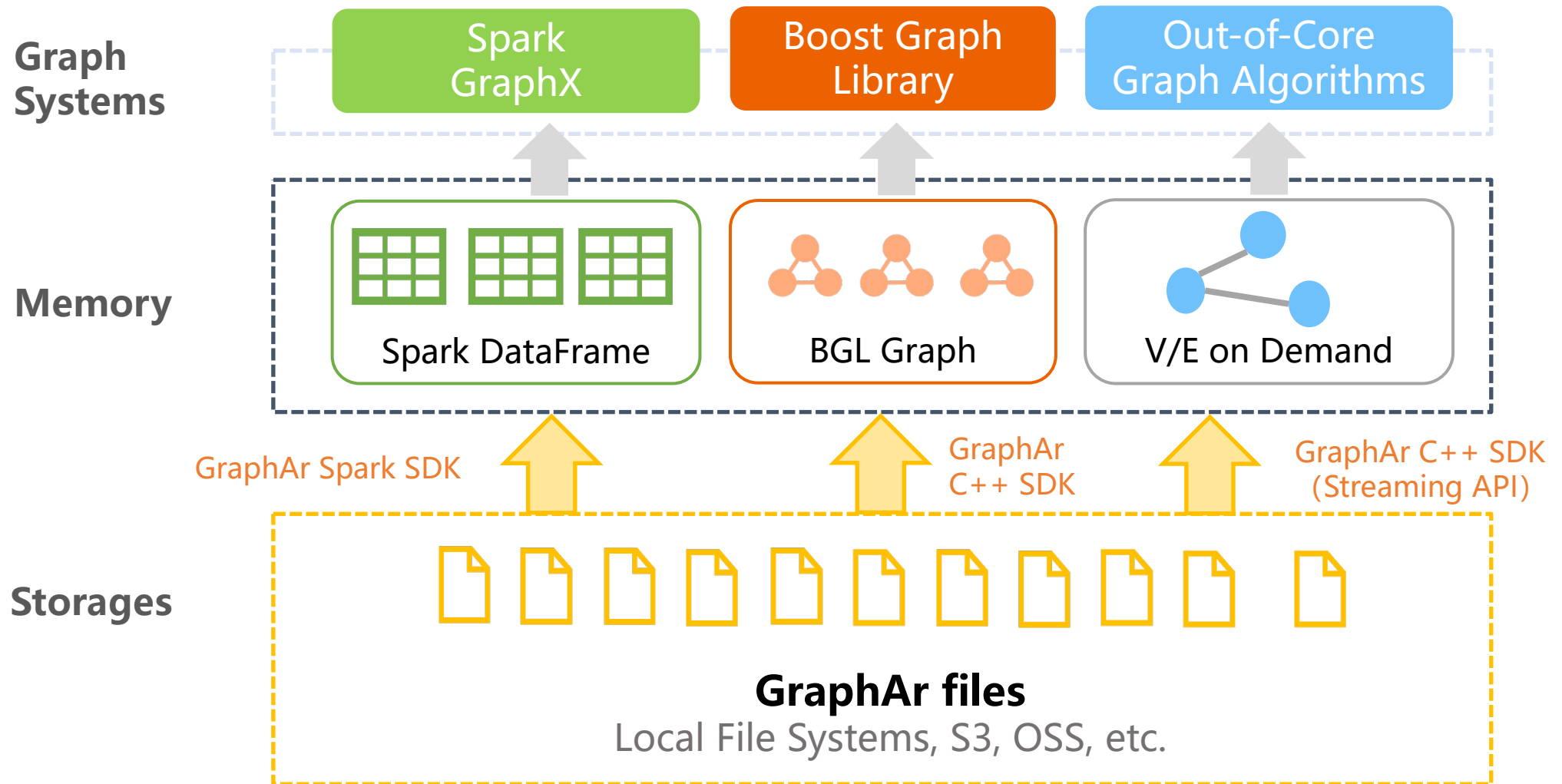
person_knows_person.edge.yml

```
1  src_label: person
2  edge_label: knows
3  dst_label: person
4  chunk_size: 1024
5  src_chunk_size: 100
6  dst_chunk_size: 100
7  directed: false
8  prefix: edge/person_knows_person/
9  adj_lists:
10 - ordered: true
11   aligned_by: src
12   prefix: ordered_by_source/
13   file_type: csv
14   property_groups:
15     - prefix: creationDate/
16       file_type: csv
17     properties:
18       - name: creationDate
19         data_type: string
20         is_primary: false
21 - ordered: true
22   aligned_by: dst
23   prefix: ordered_by_dest/
24   file_type: csv
25   property_groups:
26     - prefix: creationDate/
27       file_type: csv
28     properties:
29       - name: creationDate
30         data_type: string
31         is_primary: false
32  version: gar/v1
```

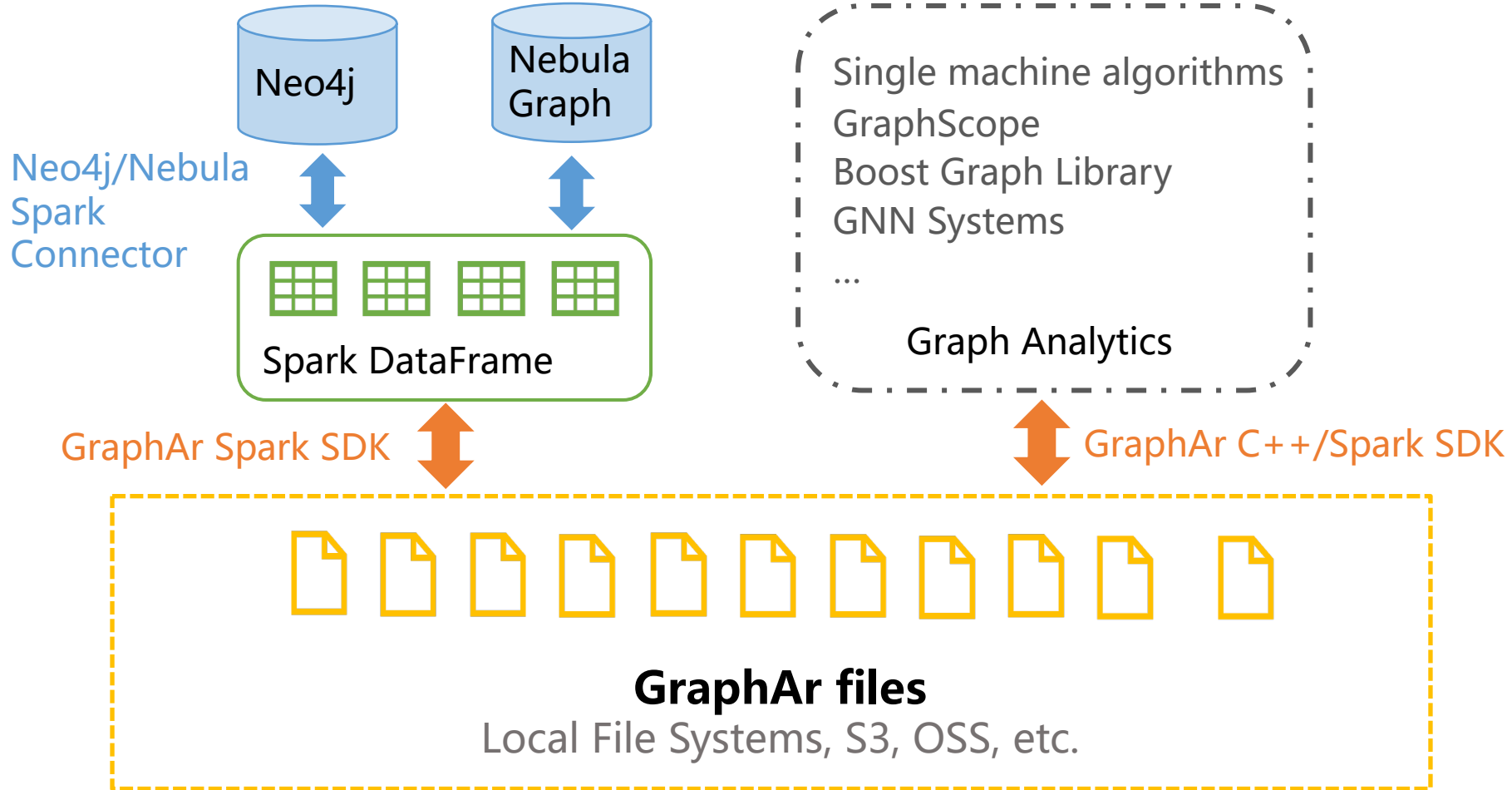
Payload files



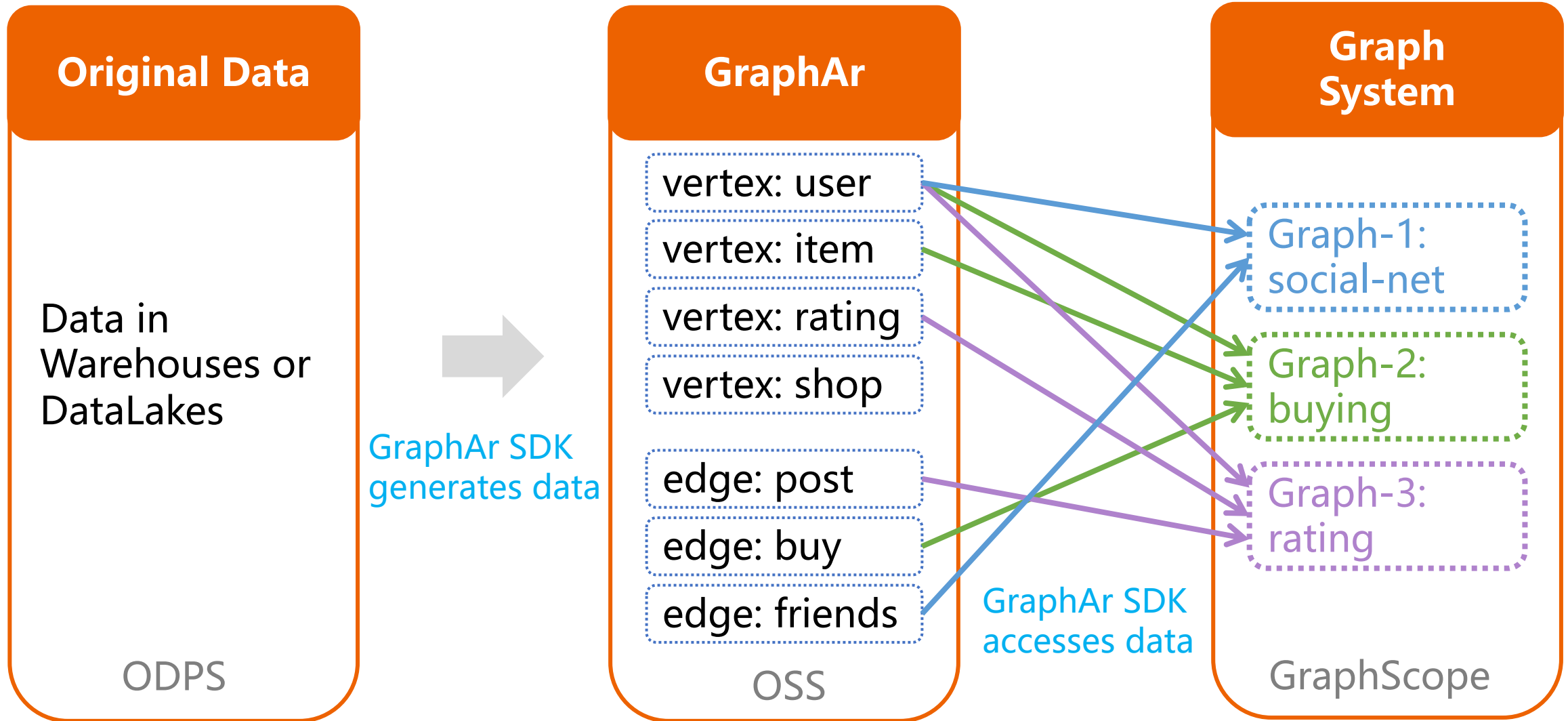
GraphAr in Action: As Graph Data Source



GraphAr in Action: As Archive of Graph DBs



GraphAr in Action: Graph Transformation



Future of GraphAr

- Support more file formats, more standard and user-defined data types.
- More graph features: RDF, time-series
- Encoding optimizations.
- Complete Spark suite to transform & create GraphAr files.
- Integrations with popular graph systems, such as Neo4j, Nebula, TuGraph, PyG ...
- Explore the use of GraphAr for data lake of graphs.

Welcome to join forces with us!



Thank you!

Learn more at <https://graphar.apache.org>

Initiated and donated by GraphScope Team from Alibaba