

Property Graphs for Industry Solution at IBM

Yinglong Xia

IBM T.J. Watson Research Center

November 14th, 2014



IBM Research



Glance of IBM Research's Efforts on Graph Technology

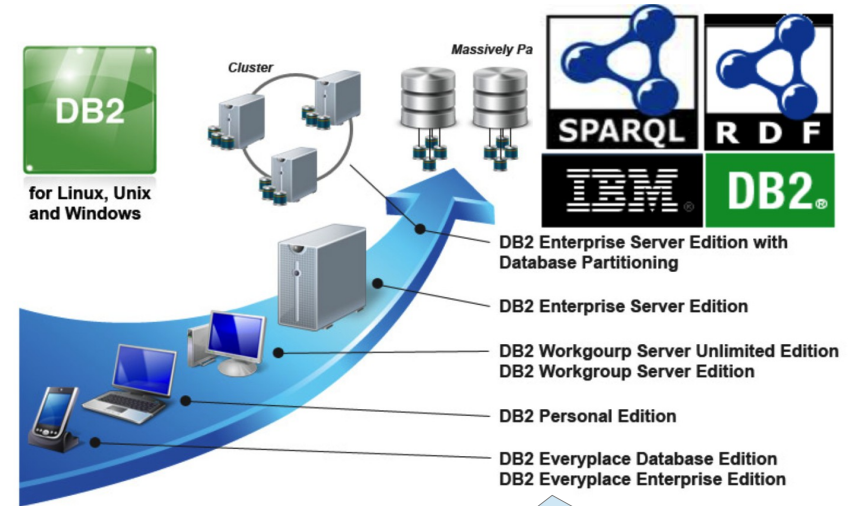


• Industry solution drives graph research

- DARPA Anomaly Detection at Multiple Scales (ADAMS) project for insider threat detecting espionage, sabotage, and internal fraud
- DARPA Social media monitoring (SMISC) for Enterprise live monitoring, thread detection, info flow analysis, and social media analysis

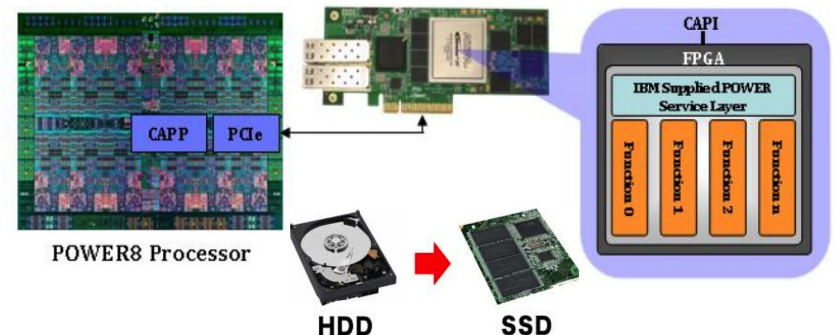
• Emerging Hardware / Accelerator

- CAPI for coherent interface
- Active store and hardware accelerators



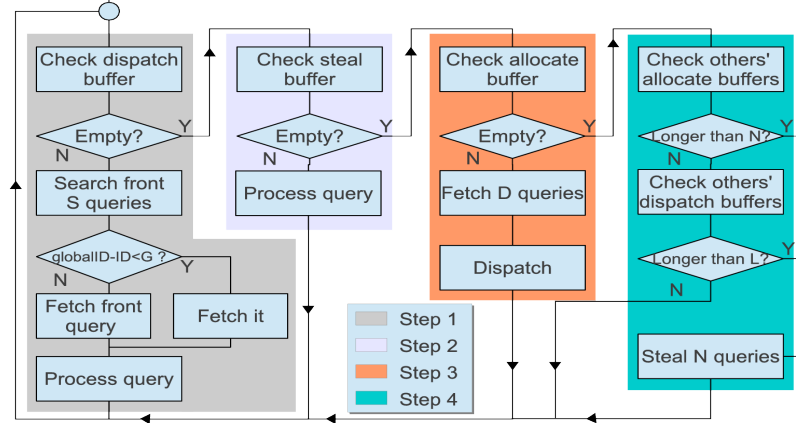
• Enhancement of IM systems

- DB2 RDF as a triple store
- DB2 Graph to support property graph query



Graph Capability in IBM System G – Native Store

User defined property class



Graph primitives

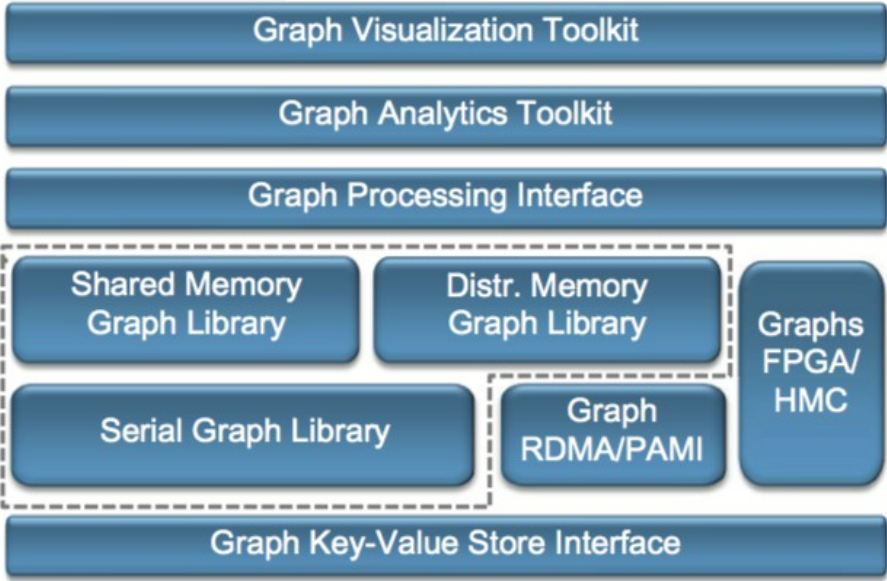
```

1 template <class VertexProperty,
2           class EdgeProperty,
3           class DIRECTEDNESS,
4           class Traits>
5 class Graph {
6     typedef ... vertex_descriptor;
7     typedef ... edge_descriptor;
8     typedef ... vertex_iterator;
9     typedef ... edge_iterator;
10    vertex_descriptor add_vertex(VertexProperty&);
11    edge_descriptor add_edge(vertex_descriptor v1,
12                            vertex_descriptor v2,
13                            edge_property&);
14    vertex_iterator find_vertex(vertex_descriptor);
15 }
16
17 typedef Graph<int, double, DIRECTED> graph1_type;
18
19 class my_vertex_property {...}
20 typedef Graph<my_vertex_property, double,
21              UNDIRECTED> user_graph2_type;
22
23 template <class G>
24 process_vertex(G& g, vertex_descriptor vid){
25     vertex_iterator vit = g.find_vertex(vid);
26     //process vertex property
27     // vit->property();
28     edge_iterator eit = vit->edges_begin();
29     for (; eit != vit->edges_end(); ++eit){
30         //process edge identified by eit
31         // eit->target(); eit->property();...
32     }
33 }
    
```

Graph traversal

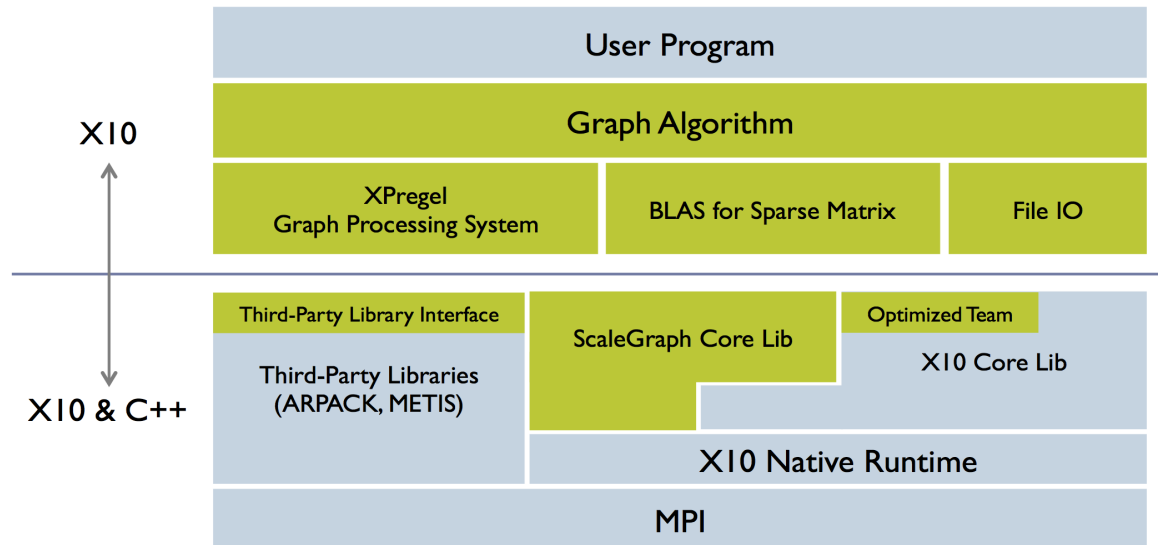
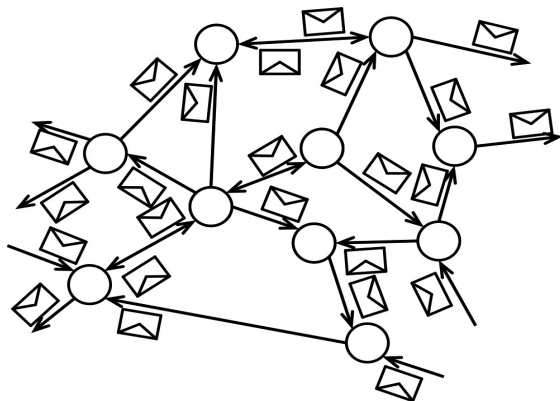
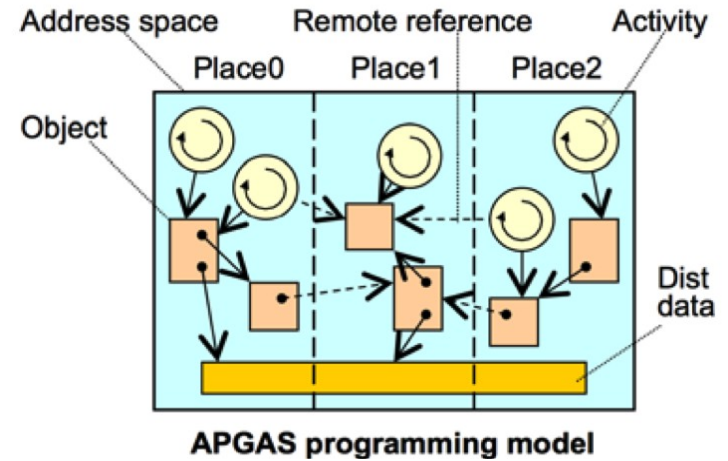
- **Native Store**
 - Performance oriented
 - Highly efficient graph access
 - Scale **up** & out

REST	SPARQL*	Gremlin
socket	TinkerPop/Blueprints	
gShell commands	JNI	
multiproperty_graph		
ibmppl_graph		
Graph KV Store		



Graph Capability in IBM System G – ScaleGraph

- IBM X10 Programming Language
 - Distributed functional language
 - Aim at productivity
 - Adopt APGAS programming model
 - Run as native code or on Java
 - Open sourced
- Graph computing model
 - XPregel model for distribution execution
 - SpVM model for graph analytics
- Performance optimization
 - Native MPI for comm.
 - Avoid serialization



DB2Graph

- Building property graph store over relational DB
 - Try to translate a gremlin query into a single SQL
 - Optimize graph update operations
 - Support batch loading of property graphs

- Shredding table
 - Without index: 46 - 230 ms per property query
 - With index: 3 - 71 ms per property query
- DB2XML
 - Without index: 100 - 2000 ms per property query
 - With index: 4 - 1000 ms per property query
- DB2JSON
 - Without index: 247 - 262 ms per property query
 - With index: 2 - 50 ms per property query

Shredding table

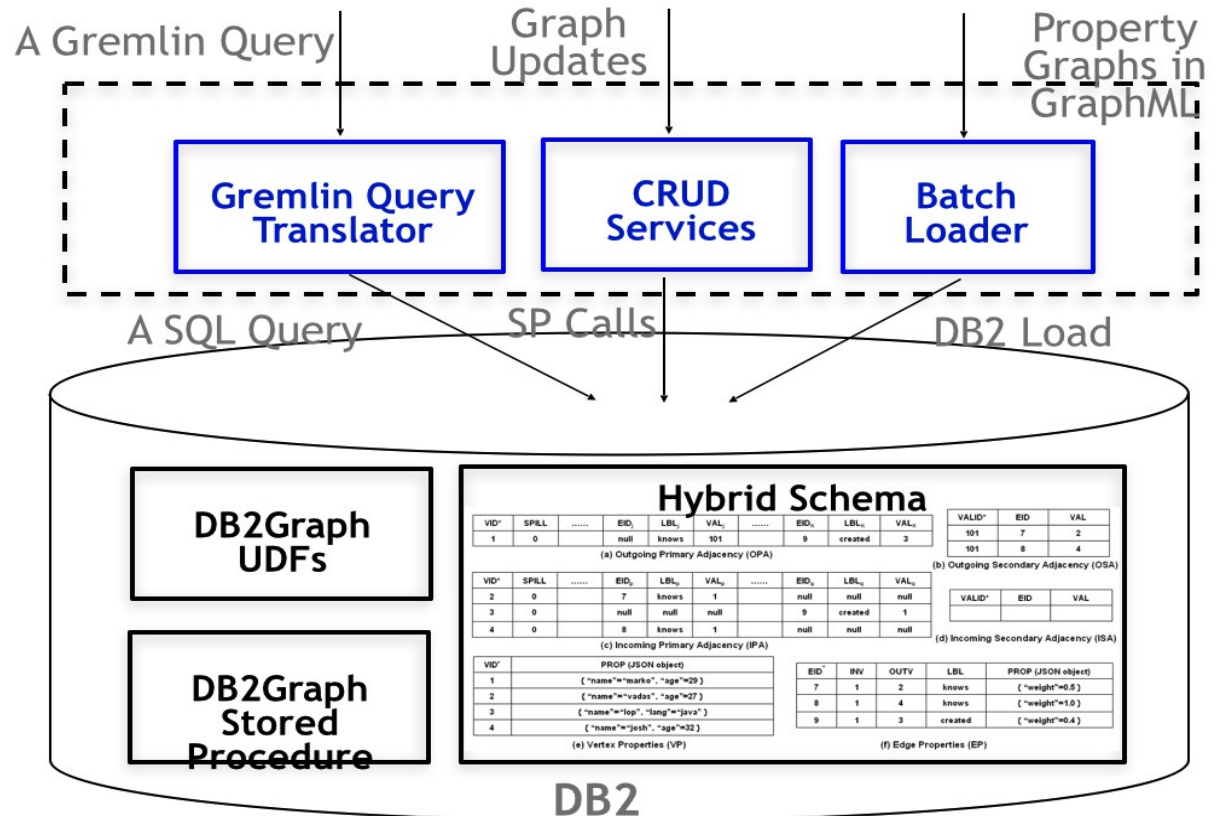
VID	PROP _i	VAL _i	PROP _j	VAL _j
1		name	marko		age	29	

DB2XML

VID	DATA (XML)
1	<pre><?xml version="1.0"?> <NodeProperty> <Property> <name>marko</name> </Property> <Property> <age>29</age> </Property> </NodeProperty></pre>

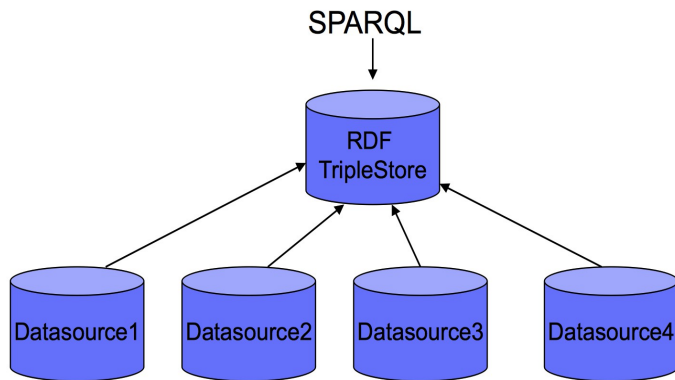
DB2JSON

VID	PROPS (JSON)
1	<pre>{ "name": "marko", "age": 29, }</pre>



DB2RDF

- In DB2 LUW 10.1 we support
 - Java API's for RDF application consumers.
 - HTTP based SPARQL query
- DB2RDF – A set of user tables within a database schema that stores the RDF data
 - Direct Primary : stores the triples and the graph they belong to indexed by subject.
 - Direct Secondary : stores the RDF objects that share the same subject and predicate within a RDF graph.
 - Reverse Primary : stores the triples and the graph they belong to indexed by object.
 - Reverse Secondary : stores the RDF subjects that share the same object and predicate within a RDF graph.



Direct Primary

Subject	Graph	predicate	value1	predicate	value4	predicate10	value10	predicate4	value40
s1		pA	vA			pB	vB	pC	lid:abc
s2				pF	vF	pB	lid:cdf		

Direct Secondary

Graph	listid	Value
	lid:abc	v1
	lid:abc	v2
	lid:cdf	v3
	lid:cdf	v4

Reverse Primary

Object	Graph	predicate	sub1	predicate	sub4	predicate10	sub10	predicate4	sub40
vA		pA	s1						
vF				pF	s2				

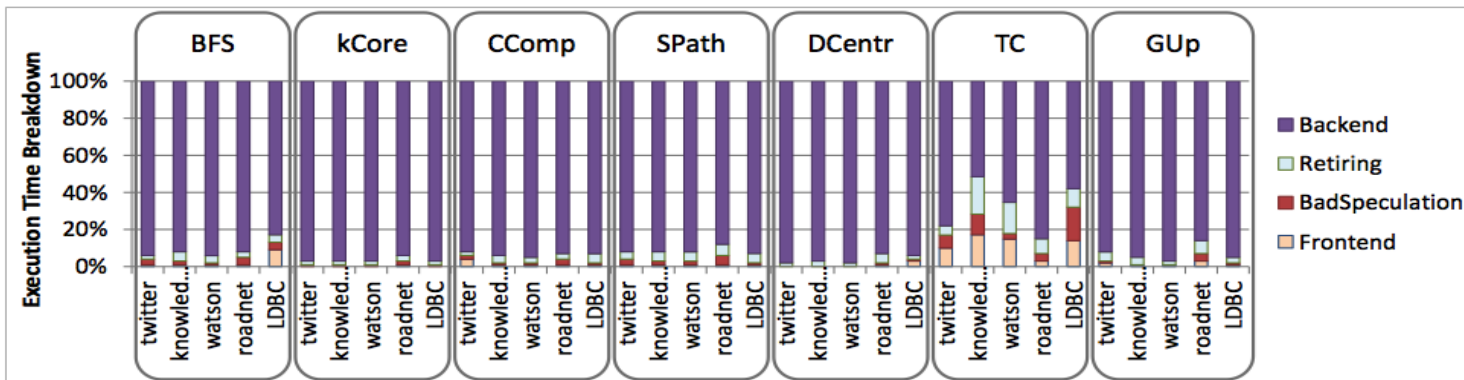
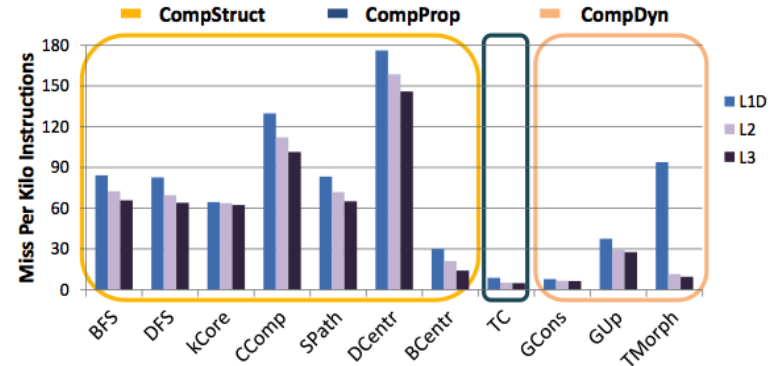
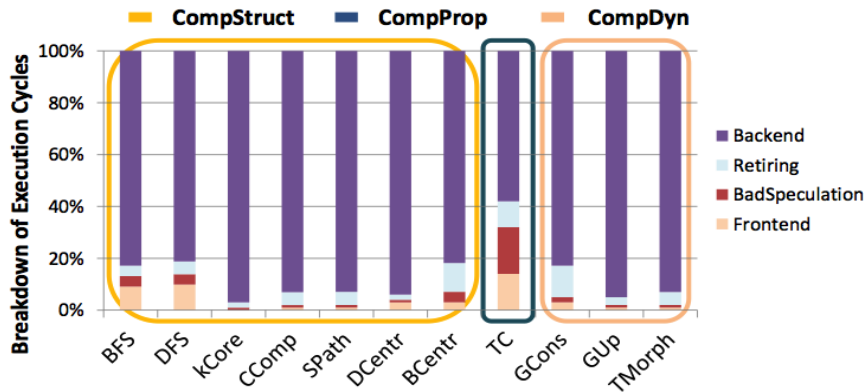
Reverse Secondary

Graph	listid	Value

Graph Benchmark for Innovations in Hardware

Experiment Data Set	Vertex #	Edge #
LDBC Synthetic Social Graph Data Set	1K	29.7K
	10K	296.1K
	100K	2887.7K
	1000K	28.82M
Twitter Graph	13M	60M
IBM Knowledge Repo	154K	1.72M
Watson Gene Graph	2M	12.2M
CA Road Network	1.9M	2.8M

- Graph for benchmarking H/W
 - Traditional processor/memory design is not friendly to graph
 - Scale **up** & out limitation
 - Demanded by novel H/W design



Performance profiling and breakdown wrt various graph algorithms`

Thank you