

PORTAL

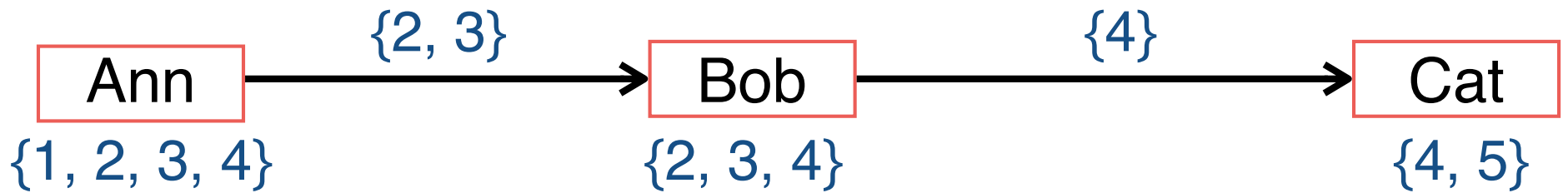
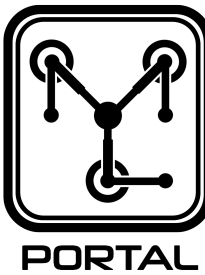
Querying Billion-Edge Evolving Property Graphs with Portal

Julia Stoyanovich
New University, USA

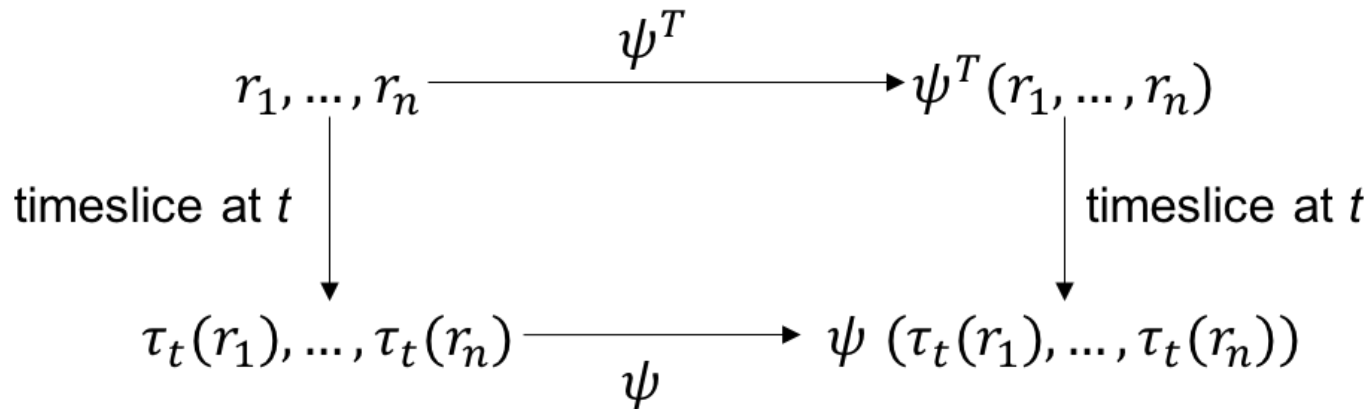
stoyanovich@nyu.edu
@stoyanoj

portalDB.github.io

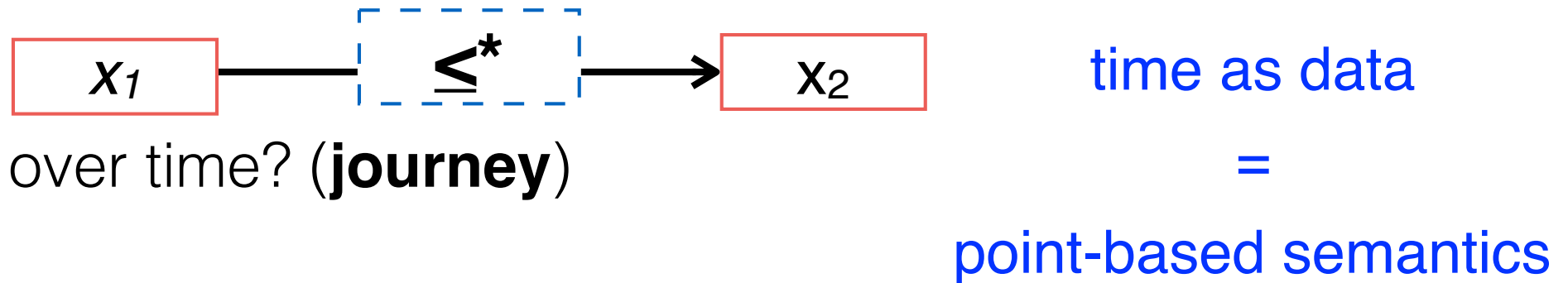
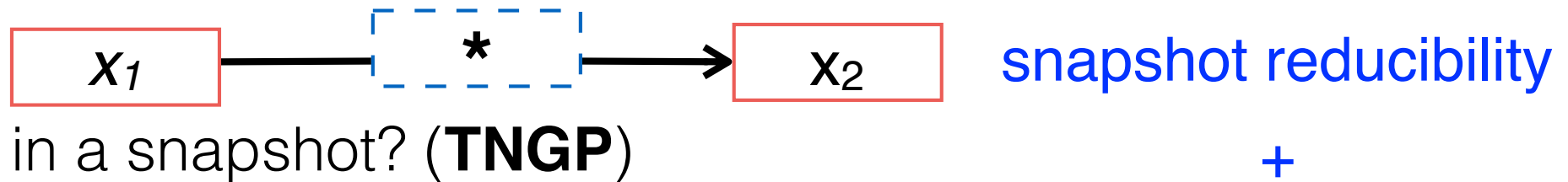
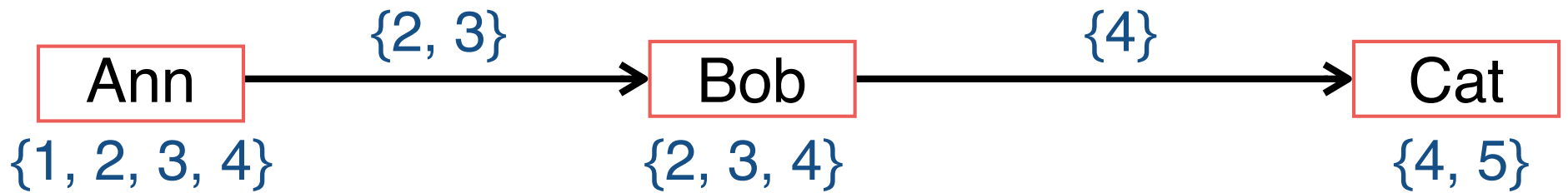
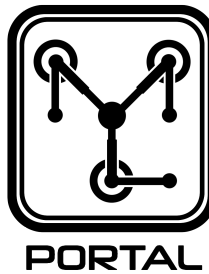
Is there a path from Ann to Cat?



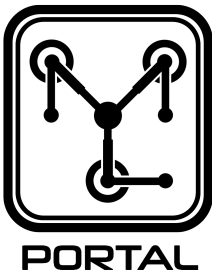
in a snapshot? (**TNGP**)



Is there a path from Ann to Cat?

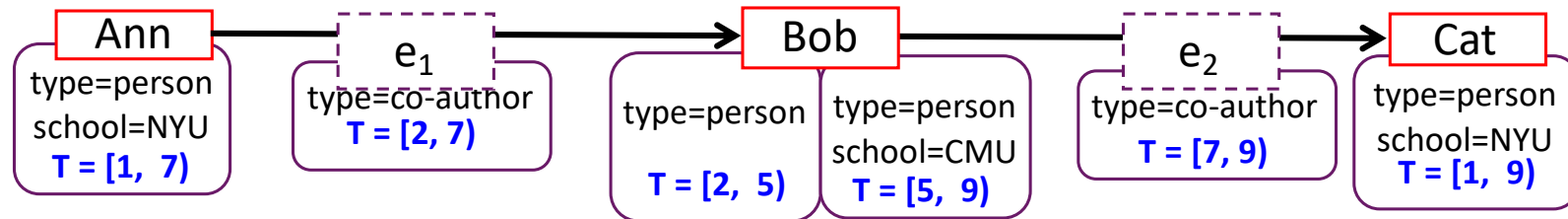
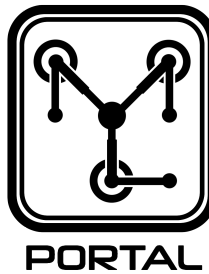


=
point-based semantics



TGraph: evolving property graph

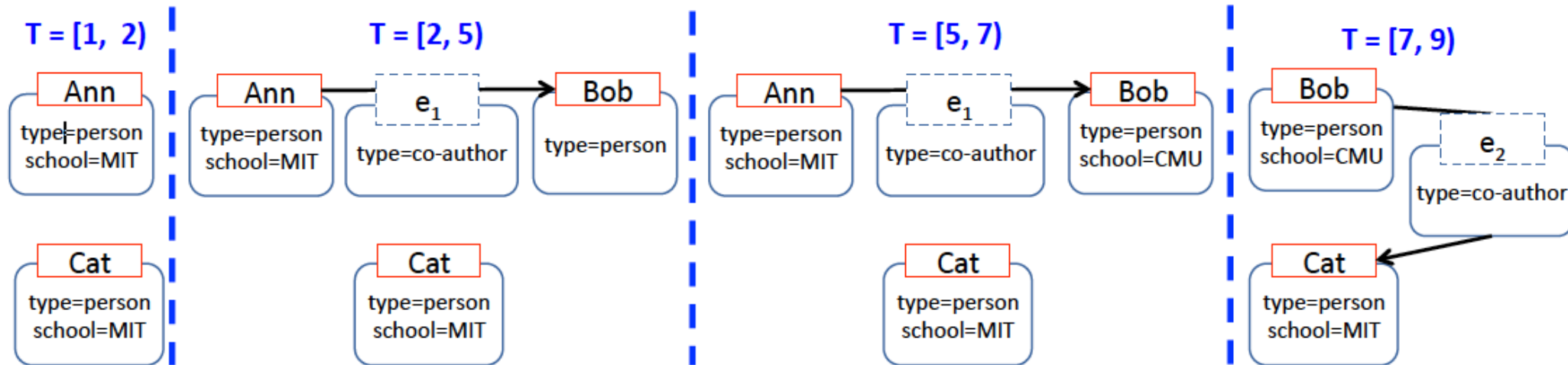
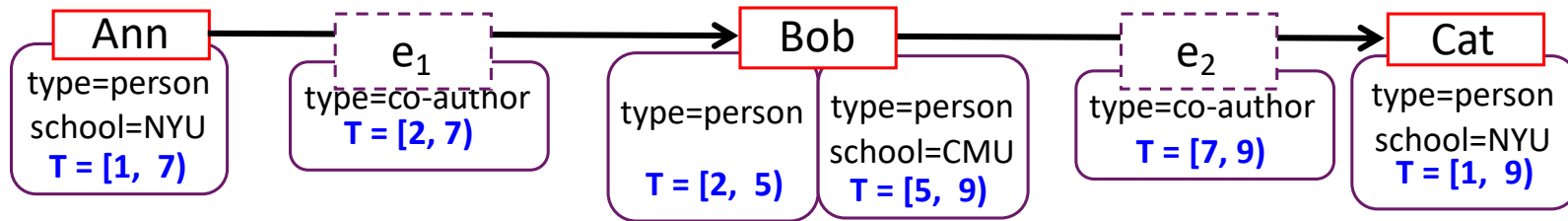
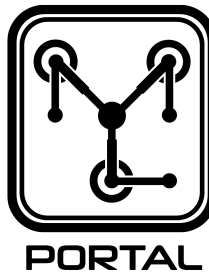
TGraph: evolving property graph



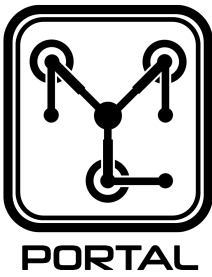
- **Definition 1.** TGraph is a six-tuple $\mathcal{G} = (V, E, L, \rho, \xi^T, \lambda^T)$, where
- V is a finite set of *nodes* (or *vertices*), E is a finite set of *edges*, $V \cap E = \emptyset$;
 - L is a finite set of property labels;
 - $\rho : E \rightarrow (V \times V)$ is a total function that maps an edge to its source and destination nodes;
 - $\xi^T : (V \cup E) \times \Omega^T \rightarrow B$ is a total function that maps a node or an edge, and time point, to a Boolean, indicating existence of that node or edge; and
 - $\lambda^T : (V \cup E) \times L \times \Omega^T \rightarrow val$ is a partial function that maps a node or an edge, a property label, and a time point, to a property value.

[Moffitt, Stoyanovich, **DBPL 2017**]

TGraph: evolving property graph

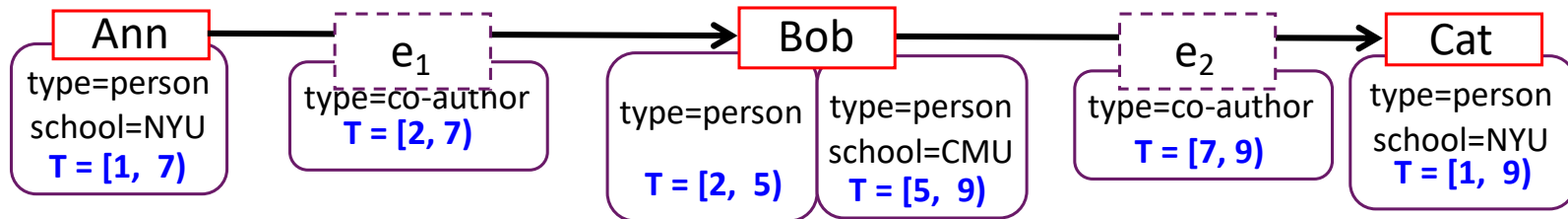
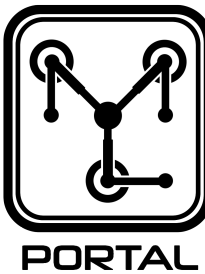


[Moffitt, Stoyanovich, **DBPL 2017**]



TGraph algebra (TGA)

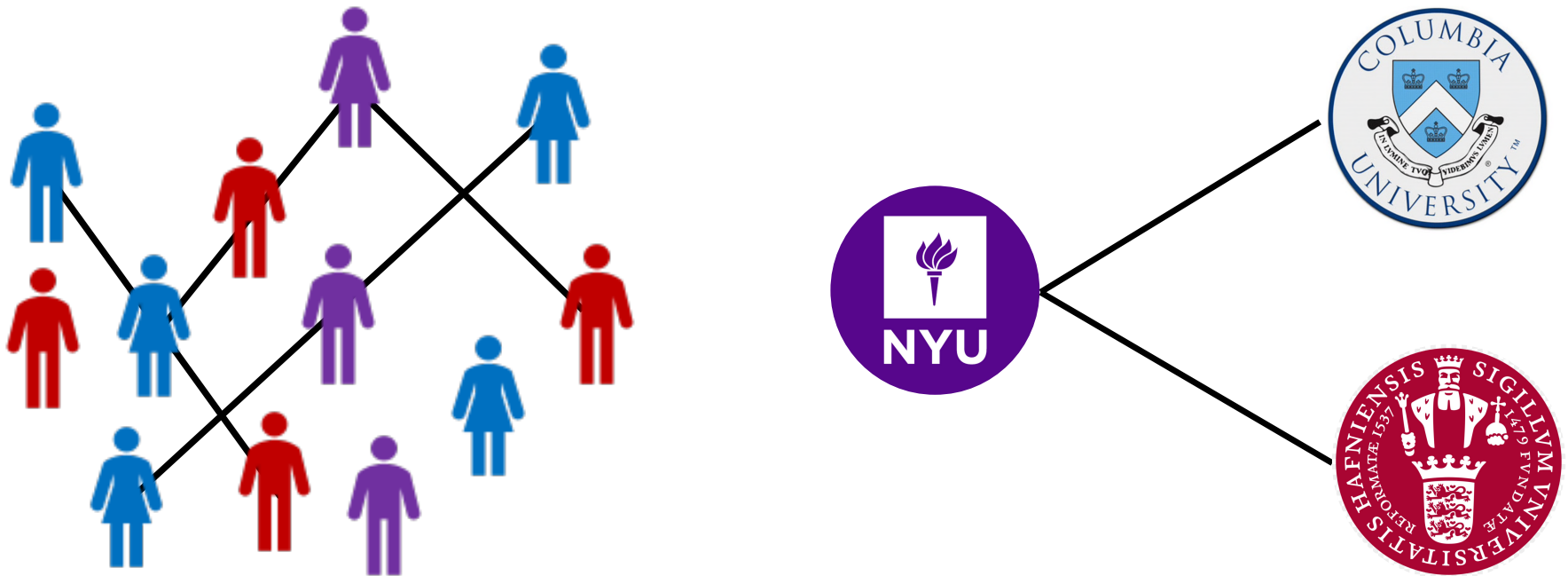
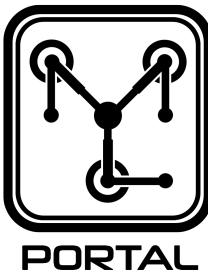
Temporal Graph Algebra (TGA)



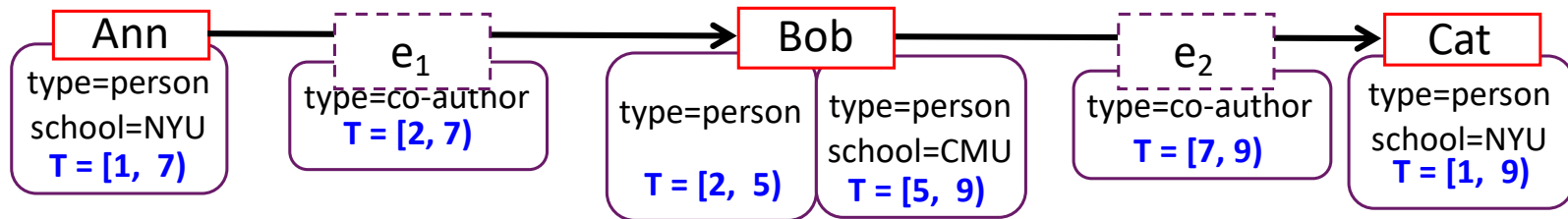
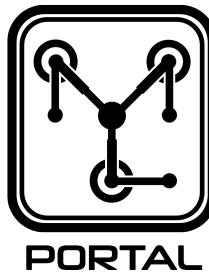
- Operators
 - temporal variants of standard graph operators: union, intersection, difference, slice, subgraph, filter, Pregel-style analytics
 - novel operator: temporal window-based zoom
- TGA is **compositional**
- Operations maintain model integrity under point-based semantics

[Moffitt, Stoyanovich, **DBPL 2017**]

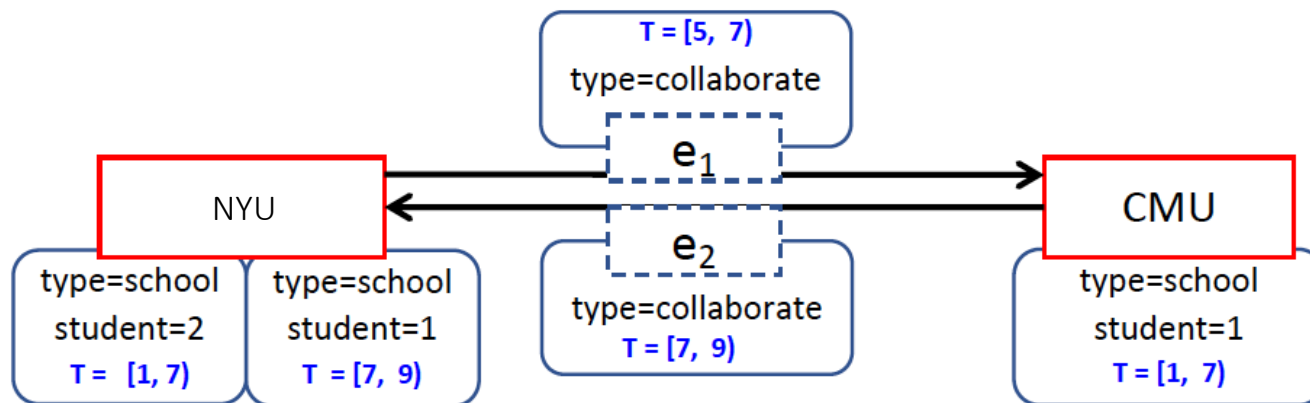
Which institutions collaborate?



Which institutions collaborate?

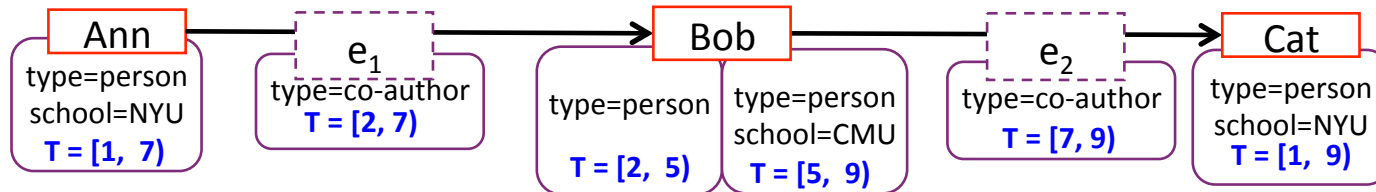
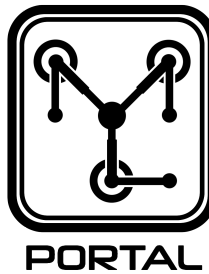


temporal attribute-based zoom



snapshot reducibility

Who is Bob's BFF?



January 2019

Sunday	Monday	Tuesday	Wednesday	Thursday	Friday	Saturday
30	31	1 New Year's Day	2	3	4	5
6	7	8	9	10	11	12
13	14	15	16	17	18	19
20	21 Martin Luther King Day	22	23	24	25	26
27	28	29	30	31	1	2

© Calendardpedia® www.calendardpedia.com

Data provided is in "without warranty"

2017

Federal Holidays 2017

Jan 1	New Year's Day	May 1	Labour Day
Jan 16	Martin Luther King Day	Jun 19	Canada Day
Jan 18	Washington's Birthday	Nov 11	Veterans Day (observed)
Feb 20	Presidents Day	Nov 23	Thanksgiving Day
May 25	Memorial Day	Dec 26	Christmas Day
Jul 4	Independence Day		

2018

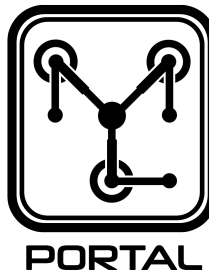
Federal Holidays 2018

Jan 1	New Year's Day	May 1	Labour Day
Jan 16	Martin Luther King Day	Jun 19	Canada Day
Jan 18	Washington's Birthday	Nov 11	Veterans Day (observed)
Feb 20	Presidents Day	Nov 23	Thanksgiving Day
May 25	Memorial Day	Dec 26	Christmas Day
Jul 4	Independence Day		

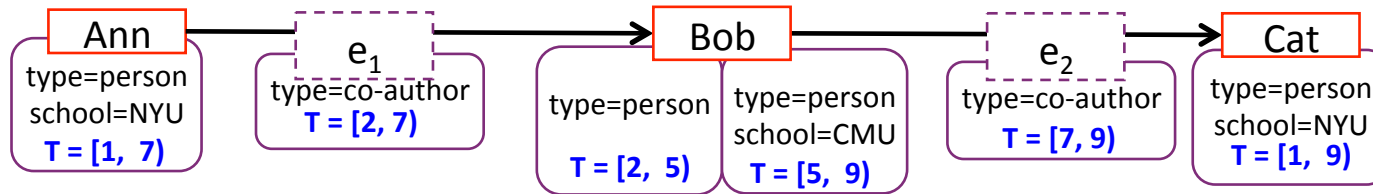
2019

Federal Holidays 2019

Jan 1	New Year's Day	May 1	Labour Day
Jan 16	Martin Luther King Day	Jun 19	Canada Day
Jan 18	Washington's Birthday	Nov 11	Veterans Day (observed)
Feb 20	Presidents Day	Nov 23	Thanksgiving Day
May 25	Memorial Day	Dec 26	Christmas Day
Jul 4	Independence Day		



Who is Bob's BFF?

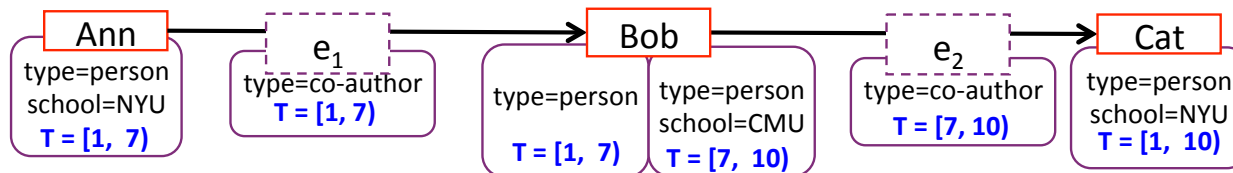


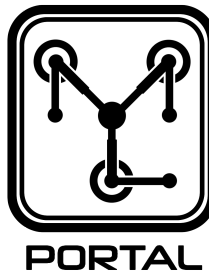
temporal window-based zoom

window	points	interval
Q1	1, 2, 3	[1,4)
Q2	4, 5, 6	[4,7)
Q3	7, 8, 9	[7,10)

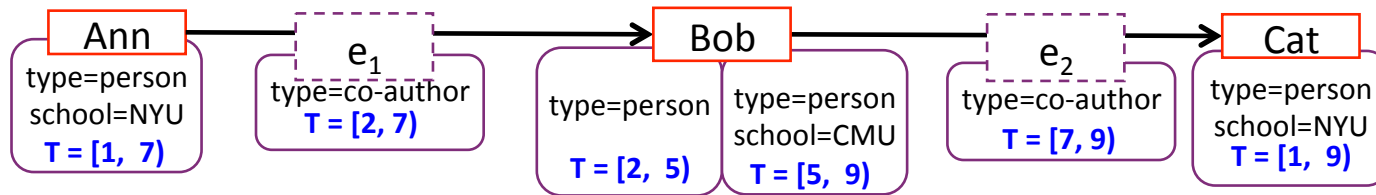
friend in *some* snapshot during a given quarter

window width=3, nodes=EXIST, edges=EXIST, node values=first, edge values=any





Who is Bob's BFF?

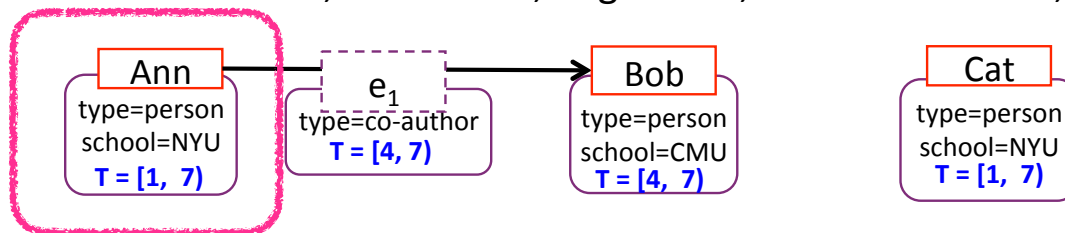


temporal window-based zoom

window	points	interval
Q1	1, 2, 3	[1,4)
Q2	4, 5, 6	[4,7)
Q3	7, 8, 9	[7,10)

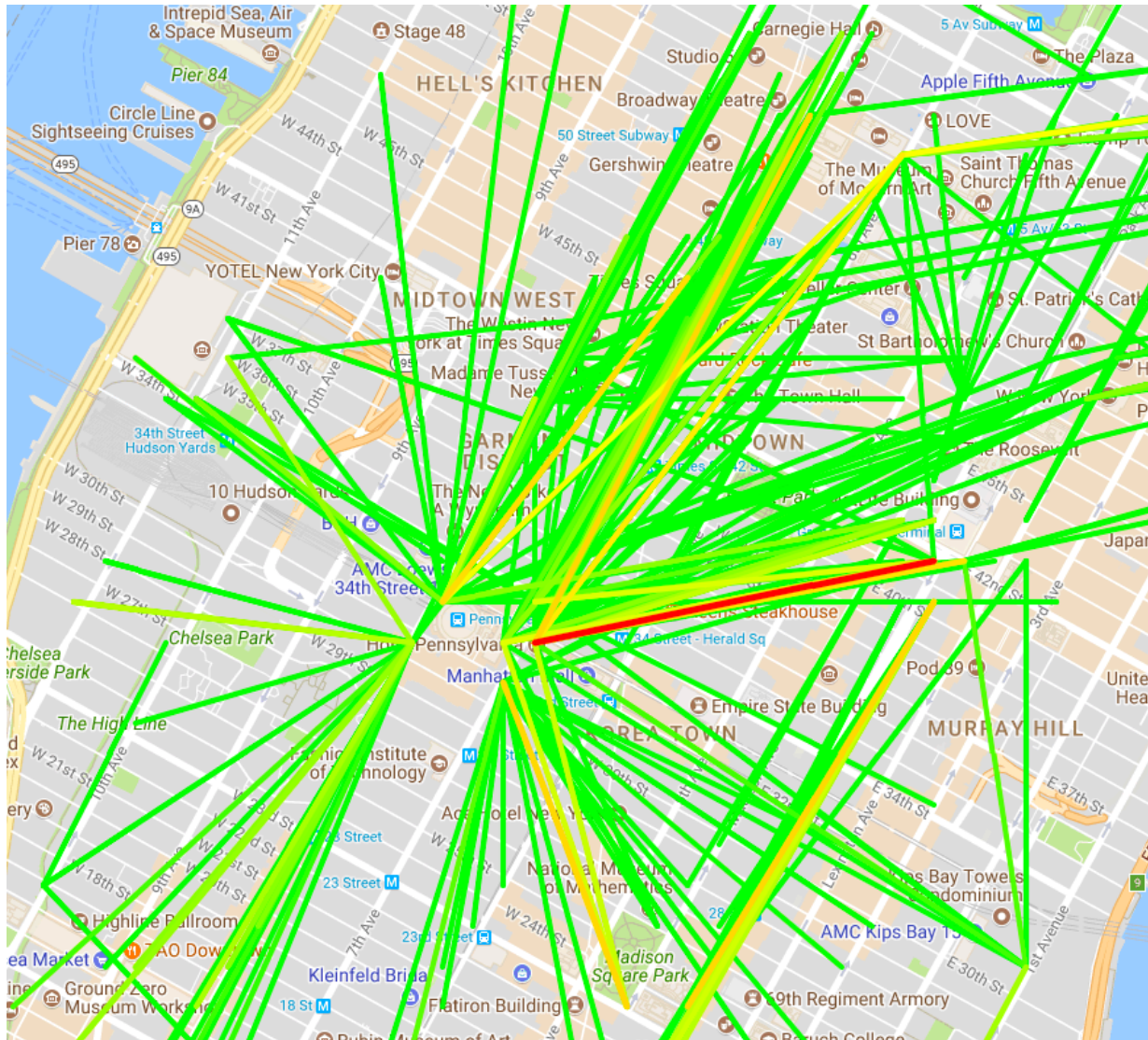
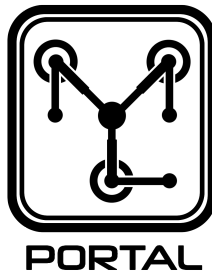
friend in *every* snapshot during a given quarter

window width=3, nodes=ALL, edges=ALL, node values=last, edge values=any



extended snapshot reducibility

Where is my bus?

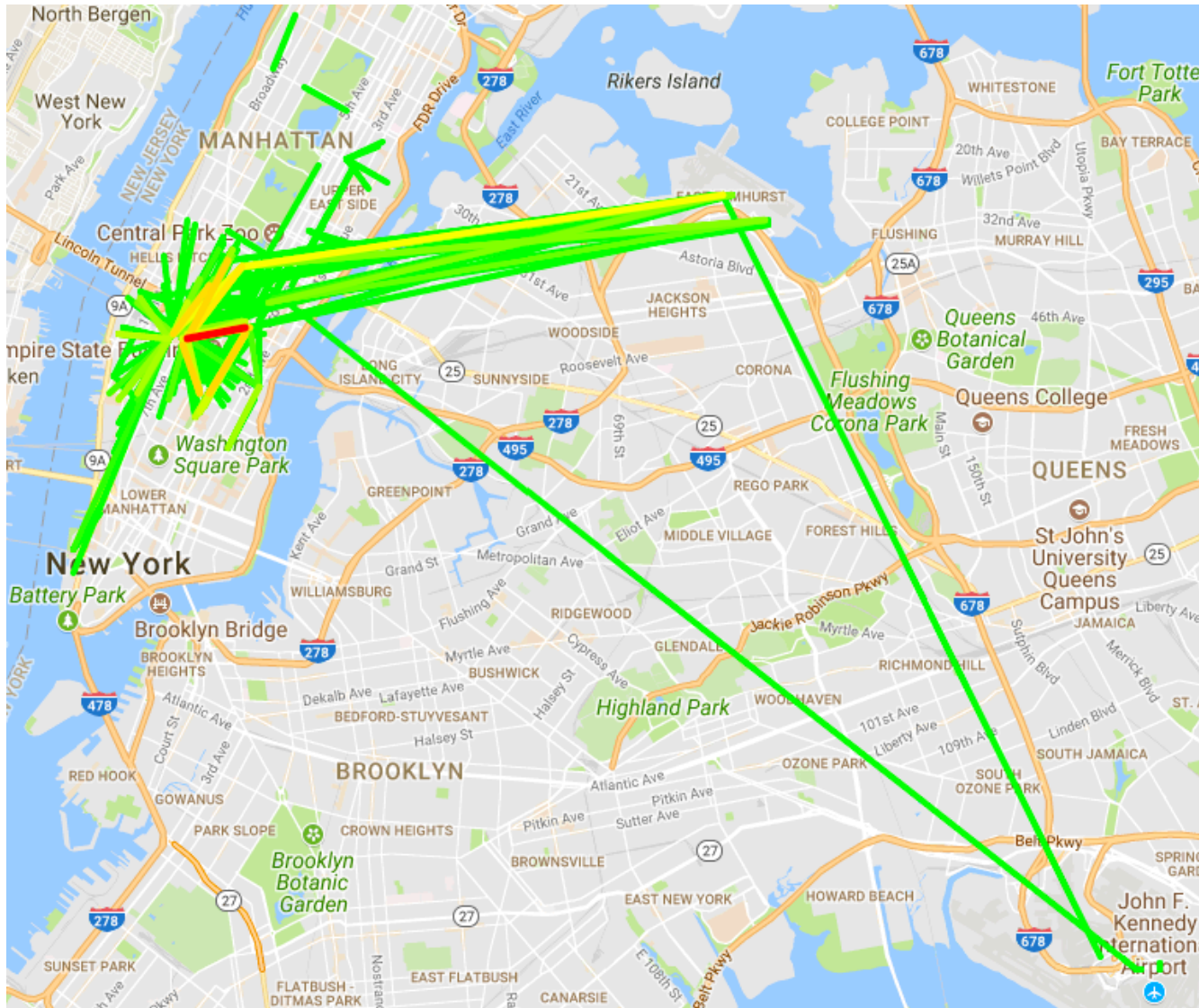
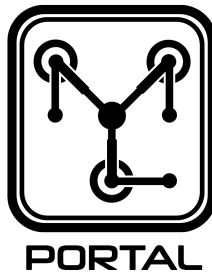


Most frequent taxi route in NYC: Penn Station to Grand Central

Why?

NYC TLC data
07/2015 - 06/2016

Where is my bus?

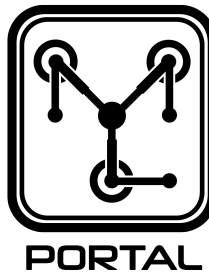


Midtown to airports:
LGA is much smaller than JFK, but gets far more taxi traffic.

Why?

NYC TLC data
07/2015 - 06/2016

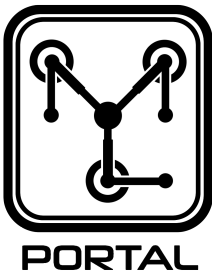
Where is my bus?



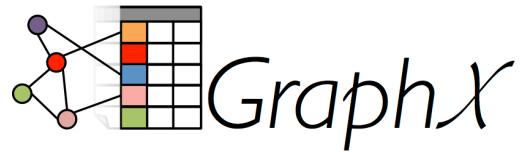
A screenshot of the NYC Taxi & Limousine Commission website. The header includes "NYC Resources", "311", and "Office of the Mayor". The main content area features a "TLC Trip Record Data" section with a map of Manhattan overlaid with yellow and green taxi trip records. A "Taxi News" sidebar on the right contains a news item titled "TLC Announces First Participant in Taxicab Leasing Pilot". The left sidebar lists various navigation options like "Home", "About TLC", and "TLC Rules and Local Laws".

- data: pick-up / drop-off time & location, fare, passenger count
- trips represented as a TGraph
 - **nodes** represent locations, with latitude / longitude coordinates as an attribute; a node exists from the time of the first incoming or outgoing trip until the time of the last trip
 - **edges** represent trips, with duration, fare etc as attributes; an edge exists for the duration of the trip

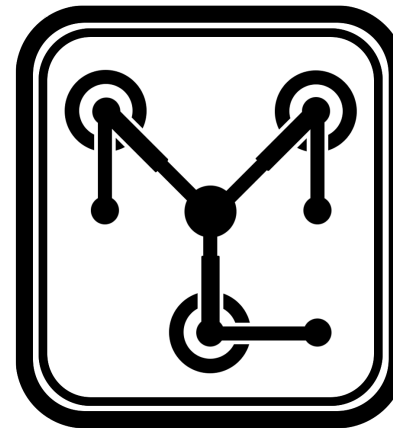
NYC TLC data
07/2015 - 06/ 2016



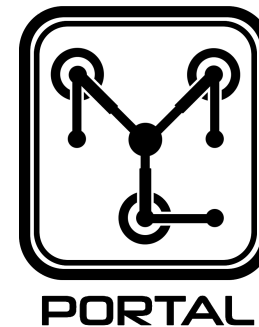
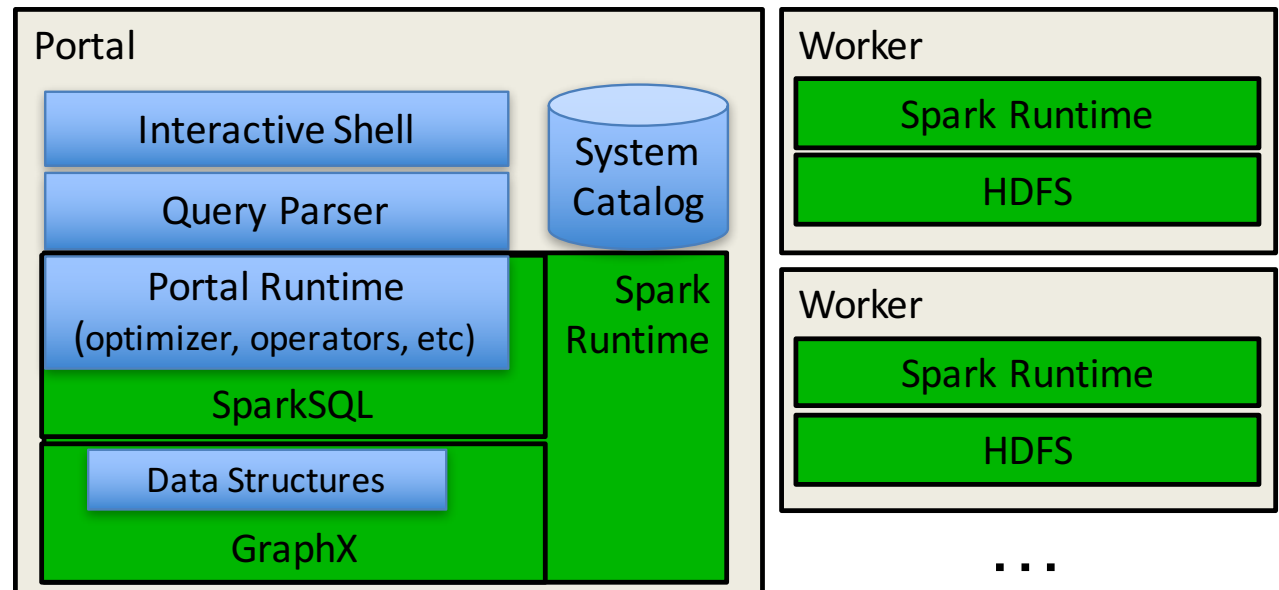
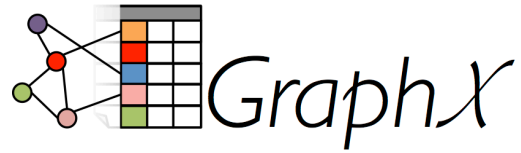
Portal: implementation



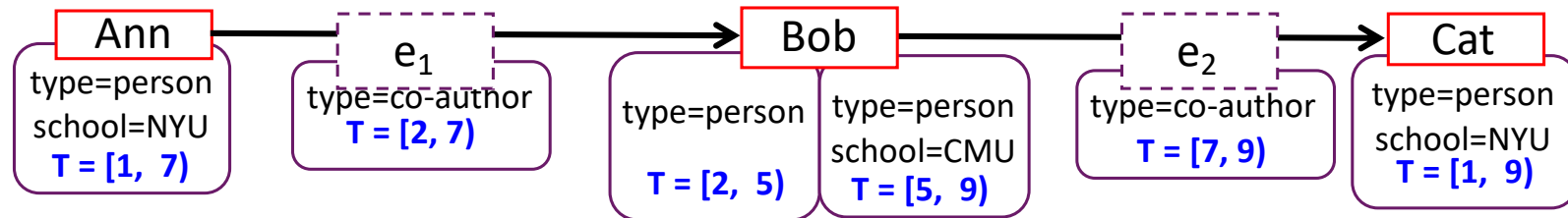
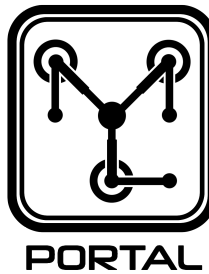
goal: principled and systematic support
for usable, scalable and extensible
analysis of evolving graphs



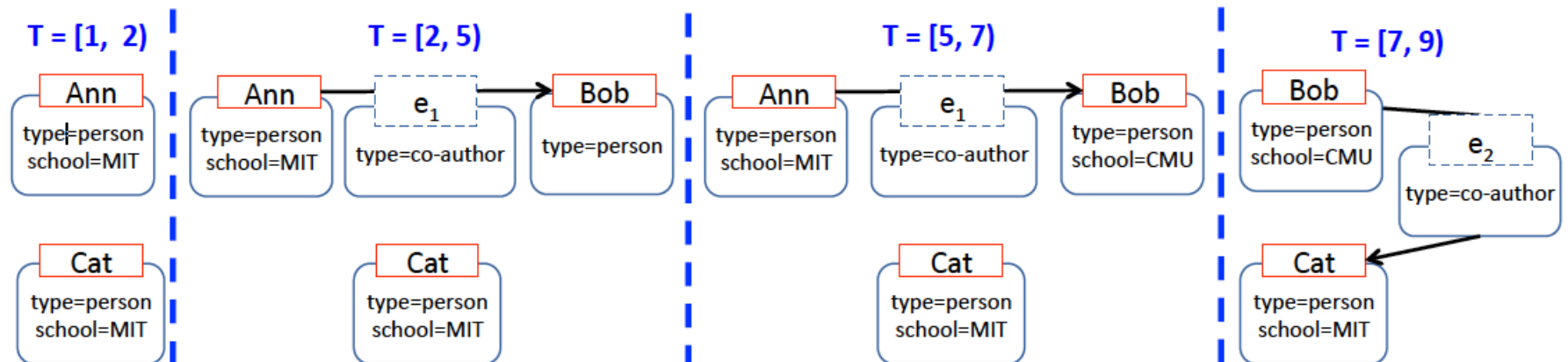
PORTAL



TGraph: in-memory representations

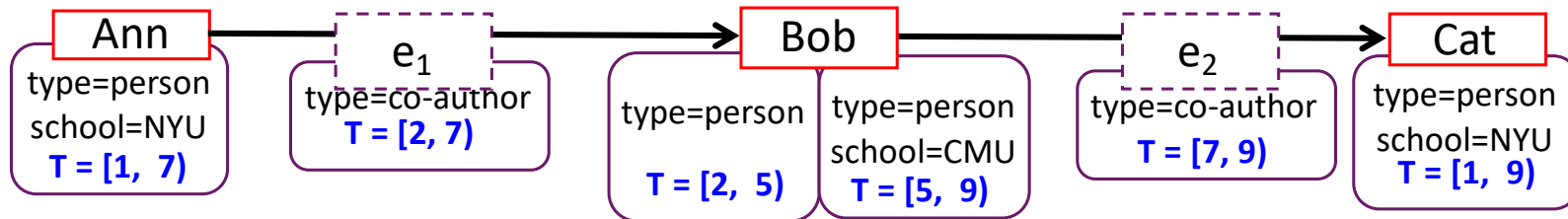
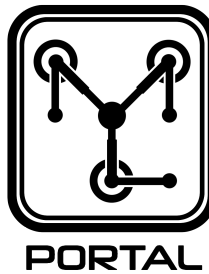


“Representative Graphs” (snapshot sequence)



[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]

TGraph: in-memory representations



“Vertex Edge” (nested relational)

Vertices

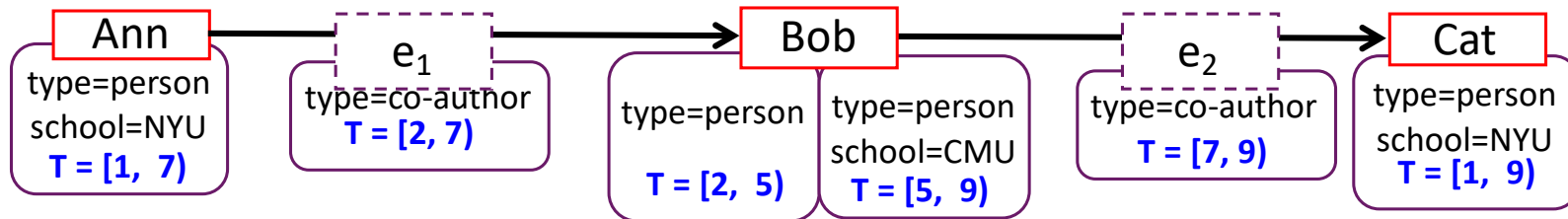
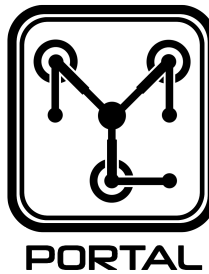
ID	Attributes	T
Ann	Type=person, School=NYU	[1,7]
Bob	Type=person	[2,5]
Bob	Type=person, School=CMU	[5,9]
Cat	Type=person, School=NYU	[1,9]

Edges

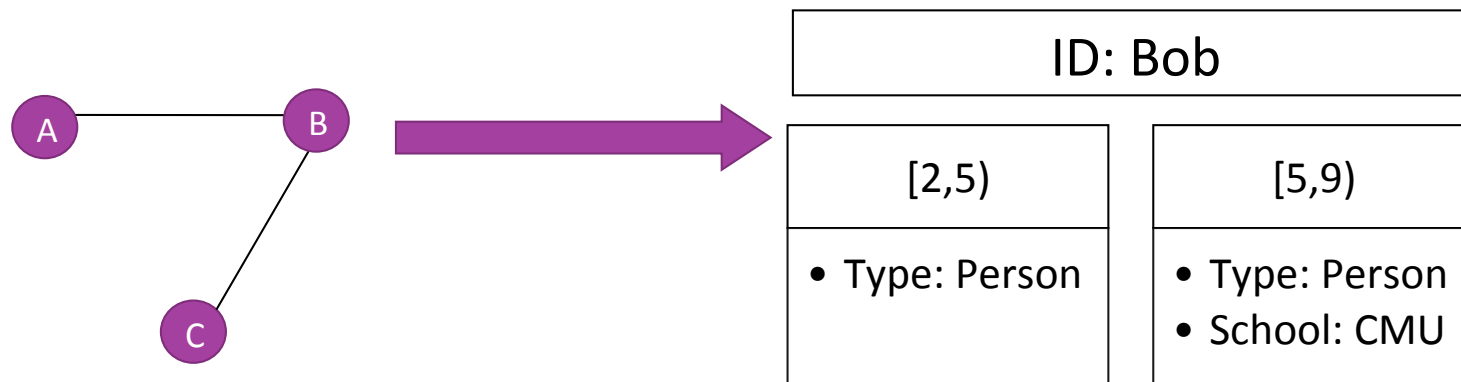
ID	V1	V2	Attributes	T
e1	Ann	Bob	Type=co-author	[2,7]
e2	Bob	Cat	Type=co-author	[7,9]

[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]

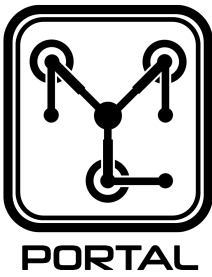
TGraph: in-memory representations



“One Graph” (GraphX graph)

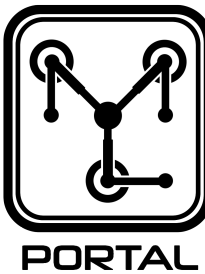


[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]



Performance highlights

Datasets



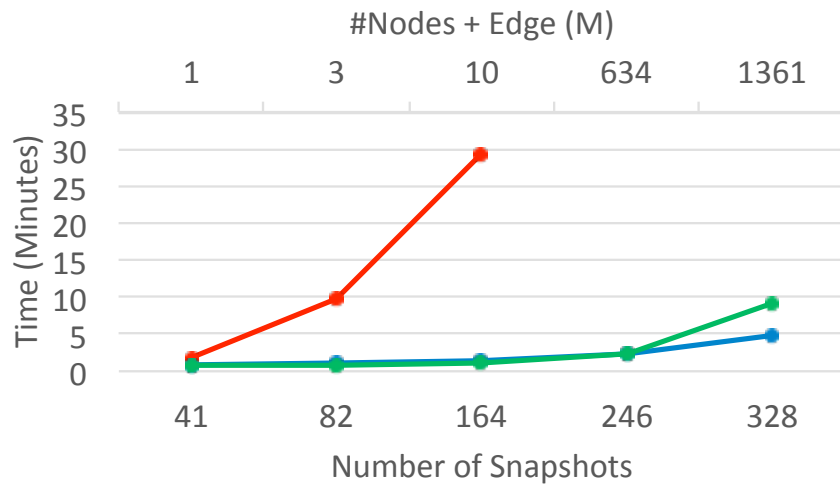
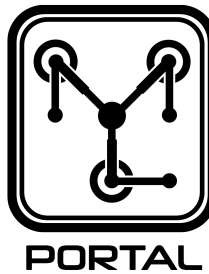
WikiTalk	SNB(LDBC Social Network Benchmark)	nGrams
Communication graph # of nodes: 2.9 M # of edges: 10.7M # of intervals: 179	Friendship graph # of nodes: 3.3 M # of edges: 202 M # of intervals: 36	Word co-occurrence graph # of nodes: 48 M # of edges: 1.32B # of intervals: 328



Cluster : 16-workers in-house cluster
Workers: 4 cores and 32 GB of RAM

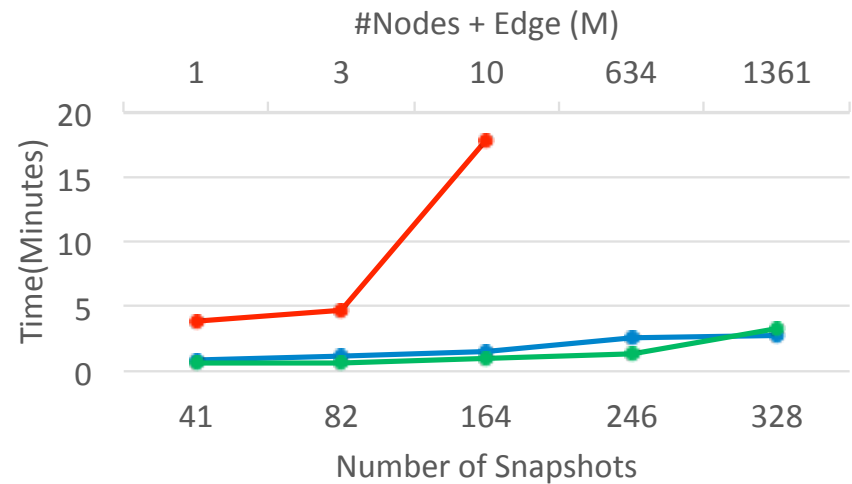
[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]

Zoom on nGrams



—●— Representative Graph —●— One Graph —●— Vertex Edge

temporal attribute-based zoom



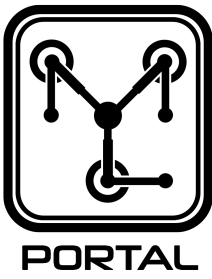
—●— Representative Graph —●— One Graph —●— Vertex Edge

temporal window-based zoom



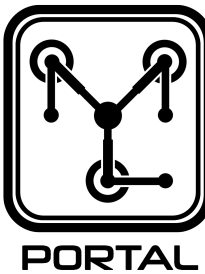
Cluster : 16-workers in-house cluster
Workers: 4 cores and 32 GB of RAM

[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]



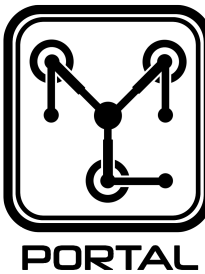
Wrapping up

Ongoing work



- **TGA / Portal**
 - Declarative language, query optimization
 - Data generation, benchmarking
 - Applications: socioeconomic studies
- **Journeys, temporal regular path queries:**
semantics, complexity of evaluation,
implementation

Take-aways



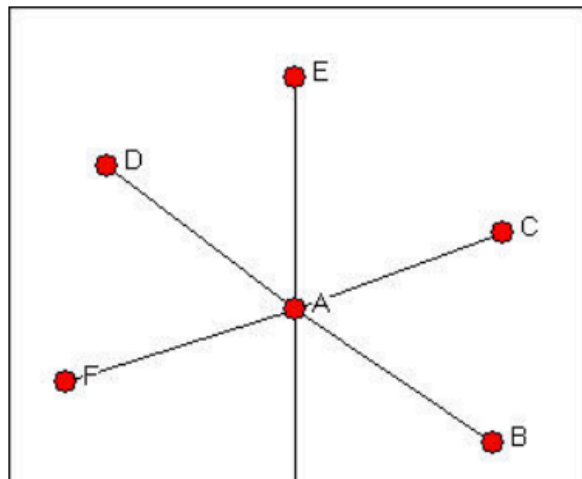
- TGraph: a logical model of property graphs with time
- TGA: a compositional temporal graph algebra under point semantics
- Portal: a library on top of Apache Spark, interoperable with SparkSQL and other libraries
- Performs well on billion-edge graphs with interesting evolution patterns
- NYC Taxi use case, working on others



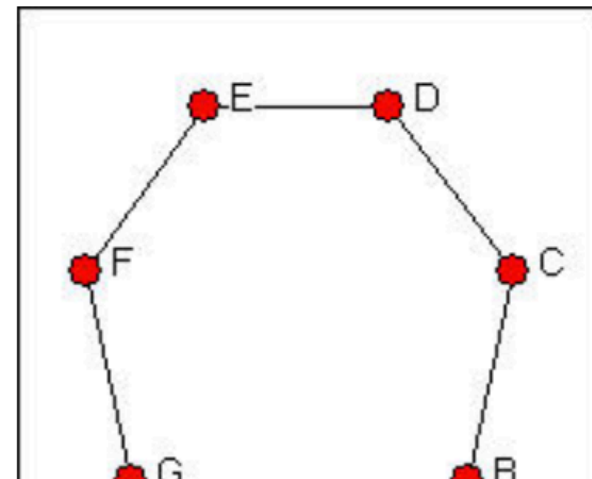
Example: Graph Centrality Over Time

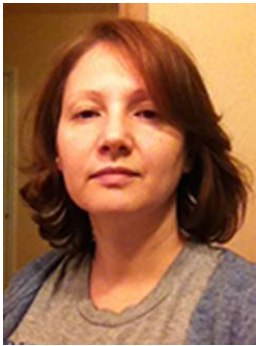
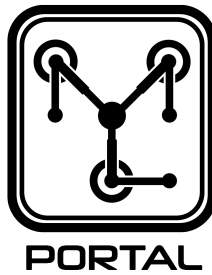
In this example we compute the degree centrality of DBLP over time. Degree centrality is a simple measure the uniformity of influence in a graph. For more information, see the definition [here](#).

The most centralized graph is a "star graph". A central node is connected to every other, and every outside node is connected only to the center.



The least centralized graph is a "circle graph", where all nodes have the same exact degree.





Vera Moffitt



Amir Aghasadeghi

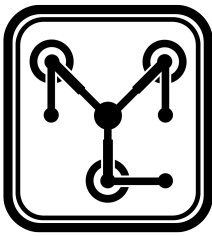


Sebastian Schelter

[CAREER] Querying
evolving graphs,
03/2018-



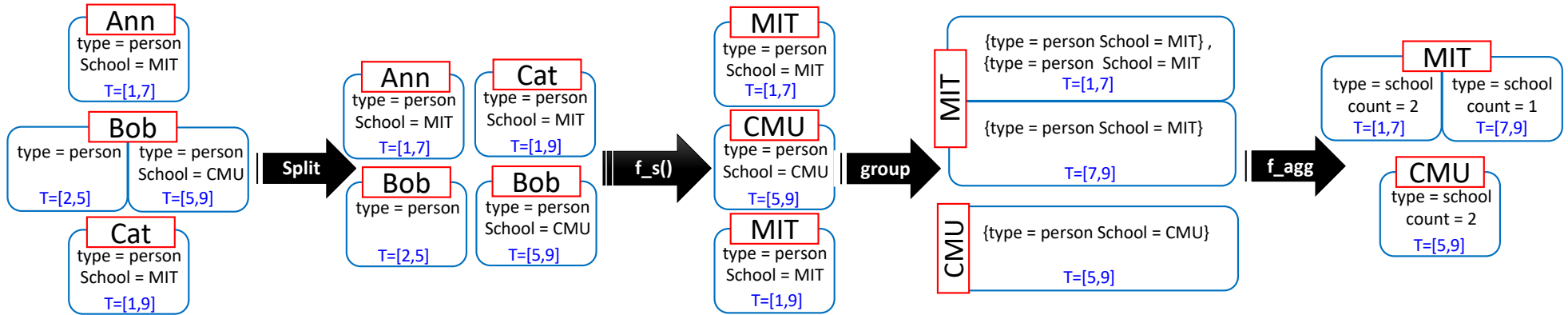
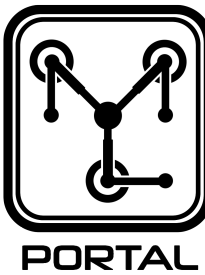
Thank you!



PORTAL

Back-up

aZoom (“One Graph”)



Algorithm 3 aZoom^T over OG

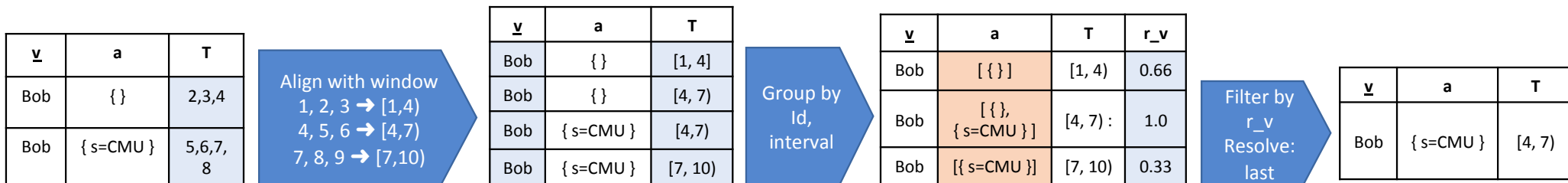
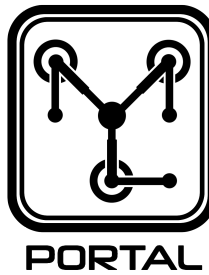
Require: Skolem function $f_s : V \Rightarrow \mathbb{N}$; Aggregation function $f_{agg} :$

$V \times V \Rightarrow V$

- 1: $V' \leftarrow V$.flatMap $\{v \Rightarrow$
 - 2: $v.history.map\{(_, attr) \Rightarrow$
 - 3: $v.copyWithIdAndAttributes(f_s(v.vid), attr)\}$
 - 4: .groupBy $\{v \Rightarrow v.vid\}$
 - 5: .reduce $\{(v_a, v_b) \Rightarrow f_{agg}(v_a, v_b)\}$
 - 6: $E' \leftarrow E$.map $\{e \Rightarrow$
 - 7: $h \leftarrow \text{recompute_history}(e)$
 - 8: $e.copyWithVidsAndHistory(f_s(e.v1.vid),$
 - 9: $f_s(e.v2.vid), h)\}$
- return** new TGraph $G(V', E')$
-

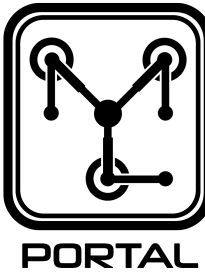
[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]

wZoom (“Vertex Edge”)



[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]

wZoom (“Vertex Edge”)



Algorithm 5 wZoom^T over VE

Require: resolve functions f_v, f_e ; quantifiers r_v, r_e

- 1: *▷ Computation of new intervals*
- 2: $I' \leftarrow I.\text{map}\{ i \Rightarrow (i, \text{computeNewInterval}(i)) \}$
- 3: *▷ Vertex aggregation for new intervals*
- 4: $V' \leftarrow V.\text{join}(I').\text{on}\{ (v, (i, n)) \Rightarrow v.n == i \}$
- 5: $.\text{map}\{ (v, (i, \text{newInterval})) \Rightarrow$
- 6: $v.\text{copyWithNewInterval}(\text{newInterval})\}$
- 7: $.\text{groupBy}\{ v \Rightarrow (v.\text{id}, v.\text{interval}) \}$
- 8: $.\text{filter}\{(i, \text{vertices}) \Rightarrow \text{match_threshold}(\text{vertices}, r_v)\}$
- 9: $.\text{reduceByKey}\{((v_a), (v_b)) \Rightarrow f_v(v_a, v_b)\}$
- 10: *▷ Edge aggregation for new intervals*
- 11: $E' \leftarrow E.\text{join}(I').\text{on}\{ (e, (i, n)) \Rightarrow e.\text{interval} == n \}$
- 12: $.\text{map}\{ (e, (i, \text{newInterval})) \Rightarrow$
- 13: $e.\text{copyWithNewInterval}(\text{newInterval})\}$
- 14: $.\text{groupBy}\{ e \Rightarrow (e.\text{id}, e.\text{interval}) \}$
- 15: $.\text{filter}\{(i, \text{edges}) \Rightarrow \text{match_threshold}(\text{edges}, r_e)\}$
- 16: $.\text{reduceByKey}\{((e_a), (e_b)) \Rightarrow f_e(e_a, e_b)\}$
- 17: **if** $r_v > r_e$ **then** *▷ Dangling edge removal*
- 18: $E'' \leftarrow E'.\text{semijoin}(V')$
- 18: $.\text{on}\{ (e, v) \Rightarrow e.\text{vid1} == v.\text{id} \text{ and } \text{in_interval}(e, v)\}$
- 19: $E''' \leftarrow E''.\text{semijoin}(V')$
- 19: $.\text{on}\{ (e, v) \Rightarrow e.\text{vid2} == v.\text{id} \text{ and } \text{in_interval}(e, v)\}$
- 20: **return** new TGraph (V', E''')

[Aghasadeghi, Moffitt, Schelter, Stoyanovich, **EDBT 2020**]