

Building Distributed Temporal Graphs From Event Streams

Benjamin A Steer
b.a.steer@qmul.ac.uk

Alessandro Di Stefano
a.distefano@qmul.ac.uk

Richard G. Clegg
richard.clegg@qmul.ac.uk

Felix Cuadrado
felix.cuadrado@qmul.ac.uk

Queen Mary University of London

ABSTRACT

Temporal graphs capture the development of relationships within data throughout time. This model would fit naturally within a streaming architecture, where new events can be inserted directly into the graph upon arrival from a data source, being compared to related entities or historical state. However, the vast majority of graph processing systems only consider traditional graph analysis on static data, with some outliers supporting either temporal analysis on similarly static data or traditional analysis on graphs updated via event streams. In this work we define a temporal graph model which can be updated via event streams and discuss the challenges of distribution and graph management. To solve these challenges, we introduce *Raphtory*, a distributed temporal graph management system which maintains the full graph history in-memory, leveraging this to insert streamed events directly into the graph model without batching and with minimal synchronisation.

PVLDB Reference Format:

Benjamin A Steer, Alessandro Di Stefano, Richard G. Clegg, Felix Cuadrado. Building Distributed Temporal Graphs From Event Streams. *PVLDB*, 11 (8): xxxx-yyyy, 2018.
DOI: <https://doi.org/TBD>

1. INTRODUCTION

Temporal graphs chronicle changes in state throughout time, unlocking a breadth of analytical possibilities. These range from expanding upon traditional graph algorithms, such as providing congestion aware GPS navigation via temporal shortest path[15], to establishing a basis for obtaining novel insights, for example investigating the long term structural changes in cryptocurrency transaction graphs [6].

Current distributed graph systems, however, focus predominantly on traditional graph processing e.g. PowerGraph [7], GraphLab [10] and GraphFrames [3], with those which do provide temporal analysis doing so in an offline batched

Permission to make digital or hard copies of all or part of this work for personal or classroom use is granted without fee provided that copies are not made or distributed for profit or commercial advantage and that copies bear this notice and the full citation on the first page. To copy otherwise, to republish, to post on servers or to redistribute to lists, requires prior specific permission and/or a fee. Articles from this volume were invited to present their results at The 44th International Conference on Very Large Data Bases, August 2018, Rio de Janeiro, Brazil.

Proceedings of the VLDB Endowment, Vol. 11, No. 8

Copyright 2018 VLDB Endowment 2150-8097/18/4.

DOI: <https://doi.org/TBD>

fashion, comparing metrics across a range of graph snapshots e.g. ImmortalGraph[11], Version Traveler [9] and GraphTau[8]. This is flawed as snapshotting reduces the granularity of temporal data to that of the snapshot window, meaning important updates may be lost. Additionally, this method does not take advantage of temporal graphs natural fit within an online/streaming environment where new events can be compared to related entities or historical state.

Kineograph[2] and Weaver[4] come close to this by streaming updates into a mutating in-memory graph model, alongside ongoing analysis. However, their approach still batches changes and analysis is performed on snapshots of the in-memory graph. Only recently have works such as Chronograph [5] began to rethink this architecture, streaming events directly into an in-memory graph.

In this paper we introduce such a temporal graph model, replacing coarse snapshots with fine-grained vertex/edge histories containing all changes to the graph structure and property values. Once defined, we discuss what a valid update to this model consists of, establishing protocols for adding, removing and updating graph entities. We explore the challenges of distributing the model across a set of machines, partitioning for high data locality and inserting updates from parallel data sources.

To solve these challenges, we present *Raphtory*, a distributed platform implementing the defined model. We then perform preliminary evaluation of *Raphtory*, investigating scalability of throughput when increasing the number of partitions and data ingesting *Graph Routers*.

2. TEMPORAL GRAPH MODEL

A temporal graph G consists of a pair $G=(V, E)$ which chronicles all transpiring changes to the graph's member vertices and edges. V is the set of all unique vertices $V=\{v_1, v_2, \dots, v_n\}$ which have existed within the graph and E is the set of all unique edges $E=\{\langle v_1, v_2 \rangle, \langle v_2, v_3 \rangle, \dots, \langle v_i, v_j \rangle\}$. An edge in this model is defined as an ordered pair of vertices $\langle v_i, v_j \rangle$, depicting directed relationships between vertices in V ; thus $\langle v_i, v_j \rangle$ is distinct from $\langle v_j, v_i \rangle$.

Each vertex in V and edge in E is assigned a chronological history $H=\{\langle t_1, created \rangle, \langle t_2, deleted \rangle, \dots, \langle t_n, created \rangle\}$ which contains all changes to the state of that entity (either *created* or *deleted*) paired with a timestamp of when the change occurred (illustrated via logical timestamps e.g. t_1, t_2, \dots, t_n). Thus, each edge and vertex exists for a set time range, or several time ranges if removed and re-added.

Each entity (vertex or edge) is additionally assigned a set of properties P , where a property is defined as a key and a value history, specifying all previous values associated with the key and the time at which the change occurred. For example, if an entity has been assigned a property key_i then $P(key_i) = \langle key_i, \{ \langle t_1, value_1 \rangle, \langle t_2, value_2 \rangle, \dots, \langle t_n, value_n \rangle \} \rangle$. These structural and property histories can then be combined to create the overall history of the graph.

To view how an equivalent non-temporal graph would appear at a given time t , a graph snapshot $G(t)$ may be generated. $G(t)$ is defined as a pair $G(t) = \langle V(t), E(t) \rangle$ where $V(t)$ and $E(t)$ represent the vertices and edges present in G at time t . For an entity to be considered present, it must be a member of its equivalent set (V or E) and its history must denote a *creation* at time t or at the closest update prior to t . Present entities will then contain a property set $P(t) = \{ \langle key_1, value_1 \rangle, \langle key_2, value_2 \rangle, \dots, \langle key_n, value_n \rangle \}$ comprising the values of all associated properties at the closest update anterior to t . $G(t)$ draws many parallels with the property graph model[12], but lacks type constraints.

2.1 Graph Updates

Updates to this temporal graph model fall into three categories: *Entity Addition* - creation of a vertex or edge; *Entity Removal* - deletion of a vertex or edge; *Entity Update* - appending entity property values. Each category has set prerequisites which must be passed before acceptance into the graph. For example, a removal is only considered valid if the entity is present, it cannot be deleted if it does not currently exist.

Entity Addition. For the addition of a given vertex v_i , it must be checked to see if $v_i \in V$. If $v_i \notin V$ it is inserted ($V' = V \cup v_i$) and a history is assigned to v_i specifying the time of its creation $H(v_i) = \{ \langle t_i, created \rangle \}$. If $v_i \in V$, its history must be inspected to see if the latest state denotes a *removal*. Given this is the case, the update is considered valid and the history is appended to specify that v_i is now present once again $H(v_i)' = H(v_i) \cup \langle t_i, created \rangle$. If v_i was not previously removed and is currently present within the graph, the update is considered invalid and is abandoned.

Edge addition necessitates similar prerequisites, requiring both the source v_i and destination v_j of an edge $\langle v_i, v_j \rangle$ to be present within V , thus avoiding hanging edges (an edge with only a source or destination). If both are present within the graph, it is checked if $\langle v_i, v_j \rangle \in E$, dictating if the edge requires insertion into E or if its history requires appending, in the same fashion as vertex addition.

Furthermore, when inserting a new entity into G the property keys it contains must be established, known as the property key-set $P_{ks} = \{ key_1, key_2, \dots, key_n \}$. Each key in P_{ks} must be assigned an initial value (iv) to be placed into the history alongside the time of entity initialisation (t_i), thus creating the property set $P = \{ \langle key_1, \{ \langle t_i, iv_1 \rangle \} \rangle, \langle key_2, \{ \langle t_i, iv_2 \rangle \} \rangle, \dots, \langle key_n, \{ \langle t_i, iv_n \rangle \} \rangle \}$.

Entity Removal. For an edge to be eligible for removal it must first be present within the graph. If so, no information is actually deleted, instead its history is appended with a *remove* state at the time at which the update occurred $H' = H \cup \langle t_i, deleted \rangle$. Vertex removal is executed in the same manner, but additionally requires the removal of all present edges within E where it is a source or destination, as these are now hanging edges; completed by appending their history with a *remove* state at the time of vertex removal.

Entity Update. To update the property value of a vertex or edge the entity must be present within G and the property key being updated must be a member of the entities key-set P_{ks} . For a given key value pair $\langle key_i, value_i \rangle$, if $key_i \in P_{ks}$ then the new value is appended into the property history along with the time the update occurred $P(key_i)' = P(key_i) \cup \langle t_n, value_i \rangle$. If $key_i \notin P_{ks}$, the update is considered invalid and is discarded.

2.2 Challenges With Distributing The Model

Correctly implementing this model in a distributed environment poses several challenges. The first issue is how to partition the graph, deciding what consists a partition and the best way to split entities spanning across them (e.g. edge-cut or vertex-cut, as described in PowerGraph[7]). However, unlike PowerGraph, partitioning a temporal graph has the additional complexity of managing tradeoffs between structural locality (proximity to neighbours) and temporal locality (proximity to an entities history) [11]. Furthermore, establishing a viable partitioning strategy for a graph built from a stream of updates is difficult as it cannot be prepartitioned and, if not actively managed, data locality will slowly degrade as more entities are added[14].

Secondly, within a distributed setting, updates are not ingested via a serial stream and may need to be sent to multiple partitions, leading to commands arriving out of order. This must be handled to ensure updates are not processed incorrectly or dropped unnecessarily. For example, if an edge add were to arrive before the addition of its source vertex, the update would be incorrectly abandoned.

Following this, updates affecting entities spanning multiple partitions must be managed/propagated efficiently. For example, if an edge-cut partitioning strategy were utilised, edges with the source and destination on different machines would require synchronisation whenever an update to the state or properties occurred. The importance of this requirement is multiplied when removing vertices which could potentially have millions of edges spanning the entire cluster, all of which would require notification of the removal.

Finally, as no updates are removed from memory, even with a large cluster of machines, eventually the memory limitations would be reached. A protocol must, therefore, be established to govern what data is retained in memory and what is offset onto more permanent offline storage.

3. RAPHTORY

To address these challenges we introduce Raptory, a system which maintains temporal graphs over a distributed set of partitions, ingesting and processing parallel updates in real time. Raptory's architecture is based on the actor model[1], with the core actor types consisting of *Graph Routers* and *Graph Partition Managers*. *Graph Routers* attach to a given input stream and convert raw data into one of the update types established in Section 2.1, forwarding this to the *Graph Partition Manager* handling the affected entity. *Graph Partition Managers* contain a partition of the overall graph, split in an edge-cut fashion. As updates arrive via the pool of *Graph Routers* the *Manager* will insert them into the histories of affected entities at the correct chronological position. This removes the need for centralised synchronisation, as commands may be executed in any given arrival order whilst resulting in the same history. An overview of this can be seen in Figure 1.

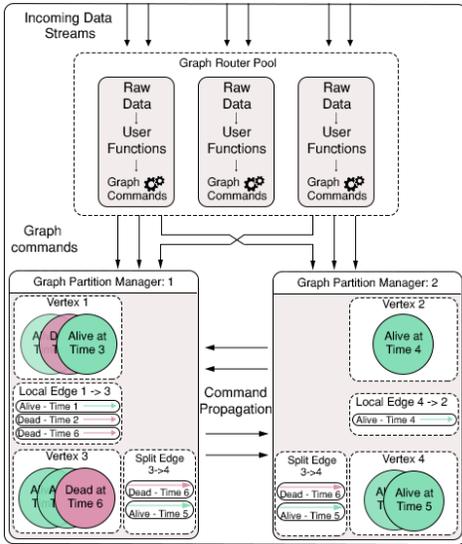


Figure 1: Raptory Architecture Overview

Raptory allows the graph to be partitioned by any given strategy. For initial testing this is a simple hash partition as it requires no state and can scale along with the number of *Routers* and *Partition Managers*. To deal with memory constraints the user may set a threshold for memory usage for each *Partition Manager*. When this threshold is met, the *Manager* will establish a cutoff point where all updates prior to this time will be transferred to offline storage. This will begin at the time of the oldest update and move forward through time at a speed proportional to the number of incoming messages. This additionally has the benefit of acting as a snapshotting mechanism to allow partitions to be recovered if the *Manager* was to crash. However, this is yet to be fully implemented and is set as future work.

3.1 Graph Router

Graph Routers are the point of ingestion for the data from which graph updates are derived. Each *Router* is attached to an input stream and converts incoming events into graph updates via user defined functions. These functions may be as complex as the user requires, ranging from fundamentals, such as what type of input is converted into a vertex or edge add, to facilitating more advanced concepts such as sliding windows or entity decay. Maintaining a larger state for advanced execution such as this may demand more resources per actor, but will never require data to be stored on disk.

When a command is generated, it is allocated a timestamp unique across all *Routers*. *Graph Partition Managers* can then use this to place the command correctly within the history of all affected entities. This timestamp is created via time fields within the raw data, under the assumption that the events were originally in the correct order at the source. However, within future work it is intended to create a method for generating unique orderings when this is not present. *Graph Routers* process events fully independently, operating a fire and forget protocol for outgoing commands, routing via the established partitioning algorithm. This allows the resources allocated for ingestion to scale dynamically according to the level of incoming data by adding or removing *Routers* from the pool as required.

3.2 Graph Partition Manager

Graph Partition Managers maintain a sub-section of the in-memory graph in the form of vertex and edge objects. These objects are organised into TrieMaps[13] (one for vertices, one for edges) and contain the entity meta data (such as its ID), its structural history and a map of its associated properties. Entity history is maintained via a red-black tree, providing fast access when reading the history and efficient insertion time for delayed or out of sequence commands. Property objects within the associated properties map mirror this structure, containing the property key and their own tree based history of previous values.

Adding Vertices. When adding a vertex, the vertex map is checked to see if an object exists for the given vertex ID. If it does, the objects history will be updated with a *Created* state and its properties will be updated with the new initial values. Note, it is not checked if the vertex is already present within the graph as a remove command may have been delayed, requiring the *Created* state to be present when it arrives. If no object within the map represents the given ID, one is instantiated, beginning the entity history and establishing the property map from a given key-set and initial values.

Adding Edges. An edge is managed primarily by the *Graph Partition Manager* storing its source vertex. If the destination vertex is stored on the same partition, the edge is considered 'local'; if it is stored in another partition the edge is considered 'split'. Edges are initialised or updated in the same fashion as a vertex, however, receipt of an edge add will also generate vertex add commands for its source and destination, avoiding possible hanging edges. These will establish placeholder vertices if the 'real' commands are yet to arrive or be ignored if they are already present. For a split edge, the *Partition Manager* will propagate the commands to the destination vertices *Manager*, requesting it to handle both this and a mirror copy of the edge. This can be seen in the edge between vertices 3 and 4 in Figure 1.

Removing Edges. When removing an edge present within the edge map, the representing objects history will be updated with a *Remove* state, again not checking if the edge is already absent as the *Created* command may arrive later. If the map does not contain an object representing this edge a placeholder entity is initialised, beginning the history with the *Remove* state. The delayed edge add may then be slotted into the history when it arrives, as well as establishing the edges property map. If the edge is split, the remove command will be propagated to the *Partition Manager* handling the destination node, as described in *Adding Edges*.

Removing Vertices. A vertex removal requires insertion of a *Remove* state into the history of the vertex and all associated edges. Unfortunately, as only existing objects can be interacted with, there is the possibility here for race conditions. Commands creating relevant edges may be delayed or received after the vertex removal and, therefore, will not contain this information within their history. For example, within Figure 1, if the command which removed vertex 3 arrived before the command which added edge 3→4, then only edge 1→3 would exist when the remove is executed. 1→3 would, therefore, be updated with the new *Removal* state, but 3→4 would miss this information.

To prevent this occurring, the *Removals* contained in adjoining vertices must be inserted into the edges history upon creation. This way, if an edge misses the execution of a

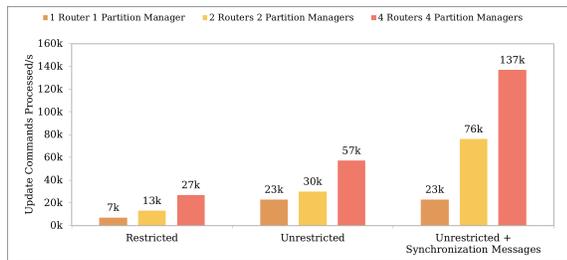


Figure 2: Initial evaluation of the maximum update throughput of *Partition Managers* in various sized clusters.

vertex removal the information is still present. For local edges, *Remove* states can simply be extracted from the history within the source and destination objects. Split edges, however, require the *Managers* storing the source and destination to exchange this information. This does generate some overhead for edge initialisation, but it is a one time occurrence and the exchange is non-blocking, inserting the removes upon arrival as if they were a normal update.

Updating properties. Properties may be updated individually or as a group. If an object representing the affected entity exists within the respective map, the new values will be placed into the history of the property object allocated for each affected key. If this update is for a split edge, it will be propagated as with addition or removal. If a property update arrives before the creation of the affected entity then a placeholder object must be utilised until it arrives.

4. EVALUATION

Raphtory actors are containerised via Docker¹ and execute independently, communicating via the Akka messaging framework². To evaluate this implementation, several secondary actor types were created, providing update generation, benchmarking and live graph analysis; these can be plugged into an established cluster without disturbing ingestion or graph maintenance. Utilising these, a stream of graph updates were generated and ingested into clusters of varying size, increasing the throughput until a maximum was reached. Three different sized clusters were tested, beginning with *Cluster*₁, a ‘singleton’ implementation with one *Graph Router* and one *Partition Manager*, and then doubling the amount of components each time; two of each (*Cluster*₂), followed by four of each (*Cluster*₄). All clusters were instantiated on servers containing an Intel Xeon E3-1284L (8 cores @ 1.8GHz) and 32GB of RAM, with one machine for each *Router/Manager* pair. The updates were then streamed in from a separate server to simulate an external event source. These events were split into 30% vertex adds, 40% edge adds, 10% vertex removals and 20% edge removals over a pool of one million unique vertex ID’s. Figure 2 displays the mean command throughput for each cluster over five iterations, split into three groups: ‘Restricted’ where each component was limited to at most one CPU core; ‘Unrestricted’ where containers could utilise all server resources; and ‘Unrestricted + Synchronization Messages’ which combined the total commands in the ‘Unrestricted’ group with

a count of the inter-manager messages used to synchronize split edges (as described in section 3.2).

When running in restricted mode *Cluster*₁ produced a maximum throughput of 7,000 msg/sec, increasing by 85% for *Cluster*₂ and a further 107% for *Cluster*₄. These are very promising results, showing that Raphtory makes efficient use of the available facilities and could work well within an shared/sparse resource environment.

When running in unrestricted mode, utilising the full resources at its disposal, *Cluster*₁ produced a maximum throughput of 23,000 msg/sec. Interestingly this only increased by 30% to 30,000 in *Cluster*₂, but soared by 90% to 57,000 in *Cluster*₄. As may be expected, upon investigation, the dampening of this initial increase was caused by the swap to a distributed environment; requiring *Routers* and *Partition Managers* in *Cluster*₂ to send inter-server updates, unlike the local communication of *Cluster*₁. Secondly, as there are no split edges in *Cluster*₁, synchronization messages are not required. On the other hand, by comparing the total number of updates processed in *Cluster*₂ against those just from the *Graph Router*, over 60% were synchronization messages. Notably though, a large proportion of these are caused by the removal of vertices with high in/out degree, generating updates for all associated edges. Removals such as this are drastic operations and it is unclear how often they appear in real datasets. This could, however, be improved by partially batching inter-manager communication and by better partitioning strategies (as discussed in Section 5). Additionally, this effect wasn’t noticed within the ‘Restricted’ tests as the clusters were limited by the amount of raw data a *Router* could parse whilst running on one core.

Overall we believe this preliminary testing demonstrates Raphtory is able to efficiently ingest and store graph events and scale with the provided resources. However, more testing will be required as the project develops, utilising real datasets, different use cases and a larger array of machines.

5. CONCLUSION

In this paper we tackle the lack of graph processing systems which perform online temporal analysis whilst ingesting new data from event streams. To address this, we firstly introduce a temporal graph model and define what may be changed in terms of graph structure and meta data. We discuss the challenges of distributing this model and provide our solution, *Raphtory*, which manages the graph history via a set of *Partition Managers*, ingesting new events into the graph via a pool of *Graph Routers*. Our preliminary experimentation shows that Raphtory scales well in both scarce and abundant resource environments, but further testing is required as the project develops.

As a continuation to this work, we first plan to implement the snapshotting/data offsetting described in Section 3, allowing for longer running tests and development of fault tolerance. Following this we intend to replace the current hash partitioning with the adaptive partition strategy described in [14], where vertices may decide to migrate to a partition with a higher number of neighbours to minimise edges spanning machines. Finally, we aim to develop ‘live graph analysis’ actors which will probe the live graph in parallel with updates, providing fast approximate results. These will be paired with ‘snapshot analysis’ actors performing the same algorithm offline to confirm the initial approximation.

¹<https://www.docker.com/>

²<https://akka.io/>

6. REFERENCES

- [1] G. A. Agha. Actors: A model of concurrent computation in distributed systems. Technical report, MIT Cambridge Artificial Intelligence Lab, 1985.
- [2] R. Cheng et al. Kineograph: taking the pulse of a fast-changing and connected world. In *ACM EuroSys*, pages 85–98, 2012.
- [3] A. Dave et al. Graphframes: an integrated API for mixing graph and relational queries. In *ACM GRADES*, 2016.
- [4] A. Dubey et al. Weaver: a high-performance, transactional graph database based on refinable timestamps. In *VLDB Endowment*, volume 9, pages 852–863, 2016.
- [5] B. Erb et al. Chronograph: A distributed processing platform for online and batch computations on event-sourced graphs. In *ACM DEBS*, pages 78–87, 2017.
- [6] E. Filtz et al. Evolution of the bitcoin address graph. In *Springer Data Science–Analytics and Applications*, pages 77–82. 2017.
- [7] J. E. Gonzalez et al. Powergraph: Distributed graph-parallel computation on natural graphs. In *USENIX OSDI*, volume 12, 2012.
- [8] A. P. Iyer et al. Time-evolving graph processing at scale. In *ACM GRADES*, 2016.
- [9] X. Ju et al. Version traveler: Fast and memory-efficient version switching in graph processing systems. In *USENIX Annual Technical Conference*, pages 523–536, 2016.
- [10] Y. Low et al. Graphlab: A new framework for parallel machine learning. In *arXiv preprint arXiv:1408.2041*, 2014.
- [11] Y. Miao et al. Immortalgraph: A system for storage and analysis of temporal graphs. In *ACM TOS*, volume 11, 2015.
- [12] Neo4J. What is a property graph. Available at <https://neo4j.com/developer/graph-database/#property-graph>, 2017. Accessed: 15-01-2018.
- [13] A. Prokopec et al. Concurrent tries with efficient non-blocking snapshots. In *ACM Sigplan Notices*, volume 47, pages 151–160, 2012.
- [14] L. M. Vaquero et al. Adaptive partitioning for large-scale dynamic graphs. In *IEEE ICDCS*, pages 144–153, 2014.
- [15] H. Wu et al. Path problems in temporal graphs. In *VLDB Endowment*, volume 7, pages 721–732, 2014.