

WHITEPAPER

# TRILLION-TRIPLES BENCHMARKING

*A Game Changer for Smart Data Management and Exploratory Analytics*

*Performance vs Previous Best Benchmark: 1.98 vs 220 hours*



## Introduction

This document describes a LUBM<sup>1</sup> benchmark at the 1 Trillion Triple<sup>2</sup> scale as performed by Cambridge Semantics Inc. on the Google Cloud Platform in October 2016.

The industry-standard benchmark executed more than *100 times faster* than any previously reported solution running the same benchmark at the same data scale. This achievement signals a paradigm shift in which a graph database rapidly moves beyond the niche “graph problem” marketplace, by providing efficient processing of diverse data at scale, to address everyday “business as usual” style analytics. This approach offers a radically more flexible and richer solution to existing solutions

The significance of this benchmark is that end users can now enjoy rapid automated database provisioning, followed by exceptionally fast data loading and from gleaning insights on vast diverse multidimensional data sets at an affordable cost.

To put these results in perspective, some examples of one trillion facts:

- ◆ 6 months of worldwide Google searches
- ◆ 100 million facts describing all the details of each of 10,000 Clinical Trial studies
- ◆ 3,333 facts describing each of the 300 million monthly active Twitter users
- ◆ 156 facts about each device connected to the internet
- ◆ 133 facts for each of the 7 billion people on earth
- ◆ 90 facts describing each of a centuries worth of Amazon orders

---

<sup>1</sup> LUBM refers to the Leigh University Benchmark <http://swat.cse.lehigh.edu/projects/lubm/>

<sup>2</sup> A triple is a simple fact consisting of a subject, a predicate and an object. See appendix 2 for more information.

## Context

The Big Data movement is characterized by rapidly increasing volumes and varieties of data and has been accompanied by much hype, driving soaring business expectations for new insights and data-derived value.

While data warehouses based on RDBMS technologies are still a predominant approach to organizing complex information for analytics and decision support they are under challenge. They are based on proven technology and governance methodologies. They offer IT a way to deliver sophisticated solutions with relatively predictable resourcing, time frames and costs. However, despite their ubiquity and the reliance placed in them by many organizations over the years, as modern data diversity and volumes grow, the cost and the inflexibility of the traditional warehouse approach is becoming increasingly impractical and in many cases unsustainable for the scale of analysis and decision support requests currently brought to IT from their businesses. The end to end time to value for RDBMS is increasingly being perceived by business as too long and too resource hungry. It results in much potential insight being left undiscovered and time to market opportunities languish.

Organizations deploying Big Data technologies are now creating data lakes containing rich information to guide their operations and decision making. These data lakes typically derive from many sources, from both inside and outside an organization and include both structured and unstructured information. The contained data generally relates to other contained data in non-obvious ways. Analytics based on these relationships can help users explore the data to discover new insights and conduct analytics to improve operational efficiency or conduct research.

The most flexible approach to link and harmonize this rich diverse data from varied sources is to use Anzo<sup>3</sup> based on the semantic knowledge graph approach that includes the data description for the pools that make up the lake.

“*The most flexible approach to link and harmonize this rich diverse data from varied sources is to use Anzo<sup>®</sup>*”

---

<sup>3</sup> Anzo<sup>®</sup> is a registered trademark of Cambridge Semantics Inc.

“  
**AnzoGraph is an in-memory graph OLAP engine that executes semantic-based SPARQL graph queries**

With a well-described “model” of the individual pools of data and how they are conceptually connected, these pools can be combined and queried in very flexible ways, even though the data wasn't originally collected with any particular integration or question in mind.

### Challenge

A data lake is becoming a critical organizational asset that grows and evolves as new diverse data sources are brought online and inter-linked. A key challenge has been how to provide the load and query performance on very large data sets drawn from the data lake in timeframes that support end users receiving an acceptable ad hoc analytics experience, against custom data set combinations and their highly individualized queries.

Software and systems benchmarks are an important aid to understanding and predicting the nature of the data lake's overall performance on a number of metrics. A standard benchmark like the LUBM is a valuable way to provide illumination of particular characteristics of the graph query engine's suitability and scalability for an organization's Anzo needs.

### Solution

AnzoGraph is an in-memory graph OLAP engine that executes semantic-based SPARQL<sup>5</sup> graph queries. AnzoGraph is based on a Massively Parallel Processing<sup>6</sup> (MPP) architecture in which compute nodes can be added as needed to address data capacity requirements and to improve query performance.

AnzoGraph can be deployed behind the enterprise firewall on dedicated enterprise servers or, as in the case of this benchmark, provisioned automatically on cloud infrastructures. The Google Cloud Platform<sup>7</sup> is one such infrastructure and provides an affordable and effective way to scale up Anzo and to provide the necessary “compute on demand” environment against even vast volumes of data as needed by end user analysts.

---

4 AnzoGraph is part of Anzo® See “Appendix 1: for additional details.

5 SPARQL is a standards based query language and protocol. <https://www.w3.org/TR/rdf-sparql-query/>

6 Massively Parallel Processing is a model for High Performance Computing (HPC)

7 <https://cloud.google.com/>

When deploying Anzo in the Cloud, the interconnect speed between the nodes is an important performance consideration. Based on this requirement, Cambridge Semantics chose to work with Google Cloud Platform to run the Lehigh University Benchmark (LUBM) for “One Trillion Facts” q.v. Appendix 2.

The result is performance more than 100 times faster than previous solutions at this scale and was operating on an industry-standard cloud platform. In other words, this workload was executed by AnzoGraph in a couple of hours, rather than more than a month’s worth of business days, as was the case with previous publications of the benchmark. Customers can use this knowledge to partially predict the behavior of their data lakes and make initial competitive analysis of the platform and software choices easier.

Performance is crucially important to unleashing the capability of organizations to meet their analytics goals in a timely, competitive manner. The balance of this document describes the benchmark results and introduces the combination of the Cambridge Semantics’ AnzoGraph and the Google Cloud Platform.

### Summary Report of Benchmark Results:

The benchmark was run October 31, 2016 on a Google Compute Platform cluster of computers. The total time for loading and querying was 1.98 hours, which was better than one hundred times faster (100X) than the previous leading vendor.

A full copy of the scripts used in the benchmark, and their output, is available to partners, customers and other qualified parties by request.

Summary Results	
587 billion triples loaded in 29 minutes and 24 seconds (total load time)	
478 billion inferred triples created in 1 hour 16 minutes and 14 seconds (total inference time)	
Query execution in 14 minutes (total query execution time)	
Total triples (load and infer): 1.065 Trillion triples	
Total time for loading and querying 1.98 hours ( <b>100x+ faster than prior benchmark</b> )	

### *System Under Test (SUT):*

- ◆ The hardware and system services were part of the Google Compute Platform (GCP) on 200 n1-highmem-32 machine type server instances.
- ◆ Each server provides 32 “vCPU”s, which correspond to 32 Intel hyper threads on 16 hardware cores. The processors are *Intel® Xeon® CPU E5-26XX @ 2.30GHz*
- ◆ The storage configured was 100GB persistent SSD drives with each node to hold the generated data.
- ◆ The nodes each had 208GB of memory available to the Ubuntu Linux 14.04 operating system and the AnzoGraph.
- ◆ The interconnect between the servers is Google’s proprietary network infrastructure. This infrastructure supports high speed communications between the database servers, communications to the Internet with VPN and live-migration for system upgrades with consistent QoS.
- ◆ AnzoGraph is the current product in full release at the time of the benchmark

The client program used to query AnzoGraph is a simple command line interface (CLI) query utility client program called “isbx”, which is included as part of the Anzo product. It ran on the lead server, along with that virtual machine’s Graph Query Engine process. Queries Q6 and Q14 were executed with compressed CSV return sets to reduce the client-server communications. The other queries used the uncompressed CSV SPARQL protocol.

## Data Generation:

The data was generated in compressed “Turtle” format, a W3C standard format for data storage and interchange.

The data was generated using the Parallel Data Generator downloaded from

<https://github.com/rvesse/lubm-uba>

Note, the data generator uses a random number generator, such that the exact triple count is sensitive to the number of servers/cores used to generate the data.

Each server contained the same number of generated data files, written to the persistent SSD drives associated with the server.

Consistent with common benchmark practice, the data generation time is not included in the benchmark timings listed above or below.

## Loading and Pruning of Duplicate Triples:

The data was loaded from the SSD drives’ compressed Turtle files in a parallel manner and automatically distributed across the cluster’s memory, indexed and pruned of duplicates, leaving just the “asserted triples”. The data was loaded as 22 logical graphs, corresponding to a data lake of 22 major sources.

The SPARQL statement “load <dir:/place/on/machine.ttl.gz> into <dst>” was executed for each of the 22 logical graphs.

LOADING SUMMARY	MEASUREMENT
Number of Triples Asserted	587,271,597,952
Total Load Time	00:29:24 hours:minutes:seconds
Effective Load	332.8 million triples per second

### Generation of Inferred Triples:

The inferred triples were generated for each of the 22 logical graphs using the AnzoGraph SPARQL- extension statement “create inferences from <src> to <dst>”. The inference (reasoning) rules and behavior are specified by the Web Ontology Language (OWL) triples present in the datasets.

INFERENCE SUMMARY	MEASUREMENT
Number of Triples Created	477,843,084,587
Total Inference Time	01:16:14 hours:minutes:seconds
Effective inference Rate	106 million triples per second
Total Triples (Load and Infer)	1.065 trillion triples

## Query Execution:

Total query execution time was slightly under 14 minutes.

Query	Seconds
Query 1	1.67
Query 2	33.3
Query 3	1.01
Query 4	3.28
Query 5	1.1
Query 6	250.27
Query 7	29.12
Query 8	37.88
Query 9	268.51
Query 10	1.56
Query 11	0.38
Query 12	11.62
Query 13	10.87
Query 14	189.18
Total	839.75

## Conclusion

At a time when organizations are struggling to eke out even incremental improvements in the time to value delivered by their data analytics programs using both the traditional RDBMS data warehouse approach as well as the more recent Apache Big Data technologies, this latest benchmark from Cambridge Semantics Inc. indicates that a third path now exists for beleaguered IT groups.

This most recent “on demand” implementation of the LUBM benchmark described in this paper, has demonstrated that a combination of the Google Cloud Platform and AnzoGraph can achieved astounding performance at data scales that dwarf most of today’s business data problems. It is exceptionally rare for any vendor to announce results that improve on previous

best efforts by two orders of magnitude. As one would expect in such a situation, the implications are profound. For the first time semantically described graph based OLAP can genuinely address mainstream data analytics challenges of all shapes and sizes. It can now deliver on fifteen years of research and development into the advantages afforded by the adoption of standards based semantic representations to the domains of data integration, data discovery and data exploration. The thorniest issue holding back the widespread adoption of semantic graph technology is now in the rearview mirror.

## Appendix 1: AnzoGraph Overview

The fundamental understanding of data relationships, reduced query complexity and scale produce a causal effect on the velocity of semantic graph databases. This is best illustrated by both the query speed and the sheer amounts of data the previously mentioned semantic graph database engine can rifle through. The combination of parallel processing and in-memory techniques maximize the discovery of relationships across a unified semantic model, enabling the parsing of billions of semantic statements each second. The significance of this fact becomes clear when considering that traditionally, issues of scale and speed (particularly in operational settings), have previously hampered graph databases. The technological advances that power contemporary query engines in semantic graph environments have addressed such concerns.

AnzoGraph is the most advanced in-memory graph OLAP database on the market and it is a critical component of the Anzo product line. AnzoGraph accommodates the wide variety of complex data that is characteristic of Big Data, enabling end users to interactively execute the rich analytics needed to uncover interesting, connected patterns in their data. Specifically it is the engine that performs the analytic queries against the facts in the data lake.

The architecture is based on the requirements of several dominant standards of the World Wide Web Consortium (W3C).

- ◆ Facts are represented in Resource Description Framework (RDF) and stored in a manner consistent with RDF
- ◆ RDF predicate URI's are specified using the W3C's Web Ontology Language (OWL) to achieve rich domain oriented data models including many relationships between many diverse entity types

- ◆ Data is queried using a full implementation of SPARQL 1.1, a query language similar to SQL, but enhanced to provide greater flexibility on potentially more complex data and richer relationships within (and across) the data sets.

These, and related standards, combine to produce the definition of a graph database system that is portable across software vendors and can easily integrate and take advantage of additional data sources from around the world.

In the context of Anzo, the SPARQL queries are dynamically generated through the configuration of ad hoc dashboarding tools guided by end user navigation of the OWL models.

AnzoGraph executes on a cluster of automatically provisioned compute servers, most typically in a cloud environment. The number of servers dynamically deployed is decided based on the data size and the response time needed for the concurrent users' queries and a fully operational cluster of arbitrary size can be achieved in under 100 seconds.

During operation, the data set is automatically distributed evenly and seamlessly across the servers. Thus, the multitude of servers and their associated CPU cores and memory constitute a single instance of the graph database.

Simply, more servers improve performance by decreasing response time and allowing more concurrent users. Adding servers allows for data volume growth and new use cases.

Anzo can support any number of AnzoGraph instances operating against part or all of the data available in the data lake. Thus, some instances operating on pools of data may be smaller, whereas instances that perform analytics across larger subsets or the entire data lake will tend to be larger.

Instances can be automatically started and shut down as needed. Lifetimes can be measured in minutes or hours for tactical instances, regardless of their size.

Therefore, data loading speed is crucially important. The "time to answer" is not just the time to query, it is often also the time to load and crucially, also the time it takes to write - or as in the case of the Anzo system, dynamically generate the query

“  
***A semantic graph is similar to a relational database but achieves greater flexibility by working with similar constructs***

Query speed is obviously crucial. Some use cases and queries must operate at human- interactive speed, others may allow a brief waiting time, and others may be less time-critical. The performance of even the least time-critical queries becomes important, though, as the organization’s needs change and the queries evolve or new queries are created and deployed to users.

### **What are Graphs?**

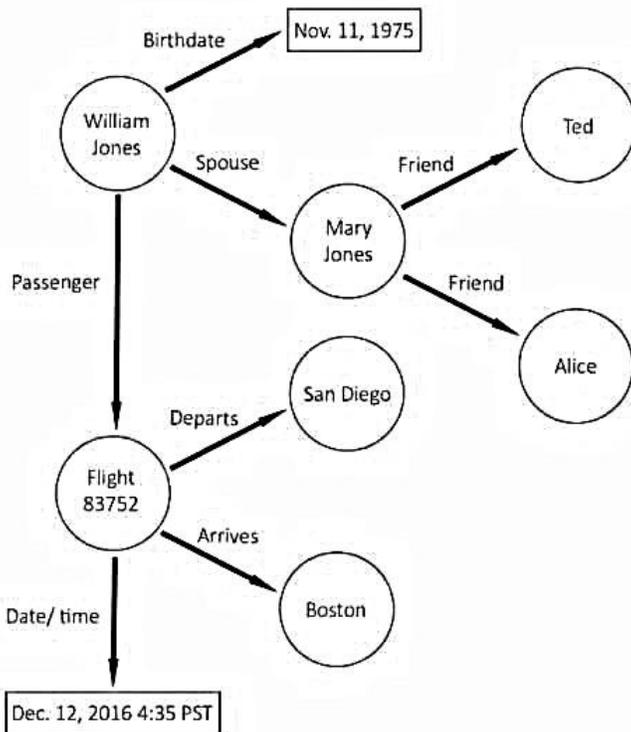
A near-endless list of books, articles, blogs and presentations exist on the subject, so this shall just provide the most cursory perspective.

The reader is advised toward the excellent primer, *“The Semantic Web and the Working Ontologist”*, industry reports from Gartner, Forrester, and other industry analysts. In addition to Semantic University and also a web search of relevant terms like RDF, SPARQL, Graph Database provide a full perspective.

Stated simply, a semantic graph is a set of interrelated facts. Each fact is a simple datum about a subject. A semantic graph is similar to a relational database but achieves greater flexibility by working with simpler constructs. As one can expect, the SPARQL query language is based on the SQL query language, but extended and tuned for the greater expressive power of graphs.

A graph is composed of some number of facts. For example, *“WilliamJones birthdate Nov/11/1975”* is a fact, typically called a triple in a graph database. Other attributes of William, such as his friends, co-workers, purchases and interests may be other triples in the database. Many things have facts, like *“Plane83572 isModel Boeing737”*. A simple data lake could contain information about planes, planned-trips, and customers. Such a database could be used for airline operational analytics, capacity planning, customer service, or customer acquisition.

Concepts may have facts, this is a very powerful extension of what has been called the *“schema”* of the data, but is now more appropriately called the *“semantic model”* or *“ontology”* of the data. For example, *“People mayHave nickName”* or *“Airplanes contain Passengers”*. These concepts describe associated physical data, and are, themselves triples. These triples can describe how data is



interrelated or provide constraints for data integrity, such as “Age mustBeLessThan 150” or “People votingAge 18”.

Further, some of the facts in a graph may not be explicitly present in the original data, but can be inferred or reasoned. For example, “WilliamJones hasSpouse MaryJones” has an inferred fact that “MaryJones hasSpouse WilliamJones”, even though the original data may only assert the former. This is because there are concept triples that specify that “hasSpouse” is a reflexive property of people.

## Appendix 2: Lehigh University Benchmark (LUBM)?

The Lehigh University Benchmark (LUBM) was created in response to the need for a standardized and systematic way of evaluating Semantic Web Knowledge Base systems. With the growth of graph technology in the enterprise space, the LUBM benchmark has been utilized by software and platform vendors to test the scalability and efficiency of graph databases as well as the reasoning capabilities of their system to support the types of questioning and analysis that need to occur.

The LUBM benchmark is based on a university ontology. “[The Benchmark’s] test data are synthetically generated instance data over that ontology; they are random and repeatable and can be scaled to an arbitrary size. It offers fourteen test queries over the data. It also provides a set of performance metrics used to evaluate the system with respect to the above mentioned requirements.”

A description of the ontology consisting of 293 triples is located at:

<http://swat.cse.lehigh.edu/onto/univ-bench.owl>

At scale factor one (SF=1), the generated data contains 103,076 triples, of which 2,531 triples are duplicated and thus the database is required to prune the duplicates to 100,545 asserted triples, including the ontology triples.

This benchmark was executed at scale factor 4,400,000 (SF=4400K). The data generator used produces 587,271,597,952 asserted triples.

At SF=1, there are 84,180 inferred triples, so the total data set is 184,725 triples.

At SF=4400K, there are 477,843,084,587 inferred triples, so the total data set is

1,065,114,682,539 triples. The fourteen queries are based on the initial SPARQL 1.0 specification and vary in both complexity and the size of the returned sets. The queries do not include the enhanced analytic capabilities defined in the subsequent SPARQL 1.1 specification, and fully supported in AnzoGraph such as aggregation, subqueries and iterative path analysis.

## About Cambridge Semantics

Cambridge Semantics Inc., The Smart Data Company®, is an enterprise analytics and data management software company. Our software, the Anzo® and AnzoGraph, allows IT departments and their business users to semantically link, analyze and manage diverse data whether internal or external, structured or unstructured, with speed, at big data scale and at the fraction of the implementation costs of using traditional approaches.

The company is based in Boston, Massachusetts

For more information visit [www.cambridgesemantics.com](http://www.cambridgesemantics.com) or follow us on [Facebook](#) [LinkedIn](#) and [Twitter](#) (@CamSemantics).

## About Google Cloud Platform

Google Cloud Platform was engineered to handle the most data-intensive work on the planet. It's the ideal environment for your business-critical applications and data, giving you the power to quickly scale and improve performance with capabilities like machine learning. And unlike other cloud providers, we're committed to building and integrating open source capabilities to preserve your control and prevent lock-in.