

1 BUSINESS INTELLIGENCE WORKLOAD

The draft BI workload was published at the GRADES-NDA workshop at SIGMOD 2018 [2]. Version 1.0 of the workload is currently being prepared.

The LDBC SNB BI workload consists of two sets of operations:

- **Read queries.** Complex read queries touching a significant portion of the data. See Section 1.4.
- **Microbatches of refresh operations.** A set of insert and delete operations, batched for a given time period (e.g. an hour, a day, etc.). See Section 1.5.

1.1 Benchmark scenario

1.2 Parameter selection

During data generation, a sequence of *substitution parameters* (??) is generated. Similarly to the Interactive workload, the parameter generation of the BI workload uses *parameter curation* [1] to ensure that the query runtimes are predictable (to some extent).

Several queries use multiple variants with different sets of input parameters. E.g. for BI 14 , 14(A) uses close countries while 14(B) uses countries that are far from each other.

In principle, query parameters are selected so that the query touches on a similar amounts of data. For queries which are only constrained by one parameter, we select ranges in the distribution where the starting node has a similar amount of neighbours. For example, if the query looks for Messages with a given Tag: (1) the Datagen computes the frequency of Messages per Tags as a factor table, (2) for each Tag, we compute its distance (absolute difference) from a given percentile of the distribution is selected (e.g. the Tag on the 75th percentile), (3) we pick the k parameters with the lowest distance.

1.3 Target metric

The performance of a system is characterized by two metrics: the geometric mean of the read query execution times and the geometric mean of the time required to load daily batches.

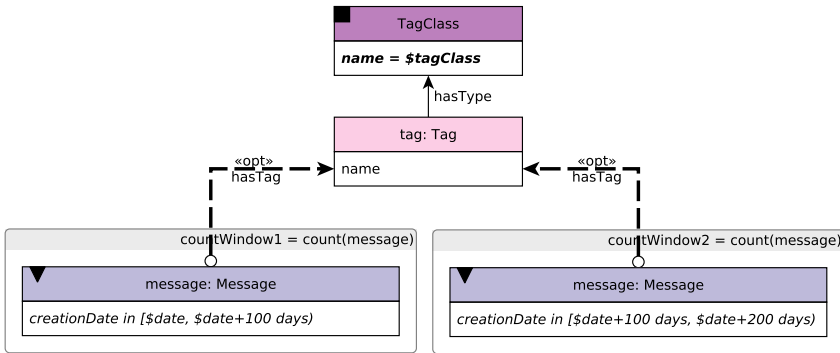
1.4 Reads

BI / read / 1

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 1				
title	Posting summary				
pattern	<div><div>message: Message</div><div>creationDate < \$datetime</div><div>length year(creationDate)</div></div>				
desc.	<p>Given a datetime, find all Messages created before that moment. Group them by a 3-level grouping:</p> <div><div><div>1. by year of creation</div><div>2. for each year, group into Message types: is Comment or not</div><div>3. for each year-type group, split into four groups based on length of their content</div></div><div><div>• 0: $0 \leq \text{length} < 40$ (short)</div><div>• 1: $40 \leq \text{length} < 80$ (one liner)</div><div>• 2: $80 \leq \text{length} < 160$ (tweet)</div><div>• 3: $160 \leq \text{length}$ (long)</div></div></div>				
params	<div><div>1</div><div>datetime</div><div>DateTime</div></div>	For later microbatches, later datetime parameters are selected keep the variance low (<0.5%)			
result	<div><div>1</div><div>year</div><div>32-bit Integer</div><div>R</div><div>year(message.creationDate)</div></div>				
	<div><div>2</div><div>isComment</div><div>Boolean</div><div>M</div><div>True for Comments, False for Posts</div></div>				
	<div><div>3</div><div>lengthCategory</div><div>32-bit Integer</div><div>C</div><div>0 for short, 1 for one-liner, 2 for tweet, 3 for long</div></div>				
	<div><div>4</div><div>messageCount</div><div>32-bit Integer</div><div>A</div><div>Total number of Messages in that group</div></div>				
	<div><div>5</div><div>averageMessageLength</div><div>32-bit Float</div><div>A</div><div>Average length of the Message content in that group</div></div>				
	<div><div>6</div><div>sumMessageLength</div><div>32-bit Integer</div><div>A</div><div>Sum of all Message content lengths</div></div>				
	<div><div>7</div><div>percentageOfMessages</div><div>32-bit Float</div><div>A</div><div>Number of Messages in group as a percentage of all messages created before the given date</div></div>				
sort	<div><div>1</div><div>year</div><div>↓</div></div>				
	<div><div>2</div><div>isComment</div><div>↑</div></div>	False < True, i.e. Posts come first and Comments second			
	<div><div>3</div><div>lengthCategory</div><div>↑</div></div>				
limit	n/a				
CPs	1.2, 3.2, 4.1, 4.2, 8.5				

BI / read / 2

BI 1	query	BI / read / 2				
BI 2	title	Tag evolution				
BI 3	pattern					
BI 4						
BI 5						
BI 6						
BI 7						
BI 8						
BI 9						
BI 10						
BI 11						
BI 12						
BI 13						
BI 14	desc.	Find the Tags under a given TagClass that were used in Messages during in the 100-day time window starting at date and compare it with the 100-day time window that follows. For the Tags and for both time windows, compute the count of Messages.				
BI 15	params	1	date	Date	Based on the creation day – TagClass – number of Messages factor table: (A) A flashmob date (B) A non-flashmob date	
BI 16		2	tagClass	Long String	For both (A) and (B), TagClasses with a similar amount of Messages are selected	
BI 17						
BI 18	result	1	tag.name	Long String	R	
BI 19		2	countWindow1	32-bit Integer	A	Occurrences of the tag during the first time window
BI 20		3	countWindow2	32-bit Integer	A	Occurrences of the tag during the second time window
		4	diff	32-bit Integer	A	Absolute difference of countWindow1 and countWindow2
	sort	1	diff	↓		
		2	tag.name	↑		
	limit	100				
	CPs	2.4, 3.1, 3.2, 4.1, 4.2, 4.3, 5.3, 6.1, 8.2, 8.5				

BI / read / 3

BI 1	query	BI / read / 3				
BI 2	title	Popular topics in a country				
BI 3	pattern					
BI 4						
BI 5						
BI 6						
BI 7						
BI 8						
BI 9						
BI 10						
BI 11						
BI 12						
BI 13						
BI 14						
BI 15	desc.	Given a TagClass and a Country, find all the Forums created in the given Country, containing at least one Message with Tags belonging directly to the given TagClass, and count the Messages by the Forum which contains them. The location of a Forum is identified by the location of the Forum’s moderator.				
BI 16	params	1	tagClass	Long String	TagClasses with a similar amount of Messages are selected	
BI 17		2	country	Long String	Big Countries are selected	
BI 18	result	1	forum.id	ID	R	
BI 19		2	forum.title	Long String	R	
BI 20		3	forum.creationDate	DateTime	R	
		4	person.id	ID	R	
		5	messageCount	32-bit Integer	A	
	sort	1	messageCount	↓		
		2	forum.id	↑		
	limit	20				
	CPs	1.1, 1.2, 1.3, 2.1, 2.2, 2.4, 3.3, 8.2				

BI / read / 4

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 4				
title	Top message creators by country				
pattern	<div><div><div>1. select top 100 forums based on memberCount in country</div><div><div>Country</div><div>name</div><div>isPartOf</div><div>City</div><div>isLocatedIn</div><div>memberCount = count(member)</div><div>member: Person</div><div>hasMember</div><div>forum: Forum</div><div>creationDate > \$date</div></div></div><div><div>2. for each country, for each of the top 100 forums (topForum1), count the Messages made by Persons who are members of any of the top 100 forums (topForum2)</div><div><div><div>topForum1: Forum</div><div>containerOf</div><div>Post</div><div>replyOf*0..</div><div>messageCount = count(message)</div><div>Message</div><div>creationDate > \$date</div><div>hasCreator</div><div>person: Person</div><div>id</div><div>firstName</div><div>lastName</div><div>creationDate</div></div><div><div>topForum2: Forum</div><div>is in top 100 forum, can be equal to topForum1</div><div>hasMember</div><div>person: Person</div><div>id</div><div>firstName</div><div>lastName</div><div>creationDate</div></div></div></div></div>				
desc.	<p>Find the most popular Forums by Country, where the popularity of a Forum is measured by the number of members that Forum has from a given Country.</p> <p>Calculate the top 100 most popular Forums. If a Forum is popular in multiple countries, it should only be calculated once with its largest membership. In case of a tie, the Forum(s) with the smaller id value(s) should be selected.</p> <p>For each member Person of the 100 most popular Forums, count the number of Messages (messageCount) they made in any of those (most popular) Forums. Also include those member Persons who have not posted any Messages (have a messageCount of 0).</p>				
params	1	date	Date	Selected from the first 30 days of the network	
result	1	person.id	ID	R	
	2	person.firstName	String	R	
	3	person.lastName	String	R	
	4	person.creationDate	DateTime	R	
	5	messageCount	32-bit Integer	A	
sort	1	messageCount	↓		
	2	person.id	↑		
limit	100				
CPs	1.2, 1.3, 2.1, 2.2, 2.3, 2.4, 3.3, 5.3, 6.1, 8.2, 8.4				

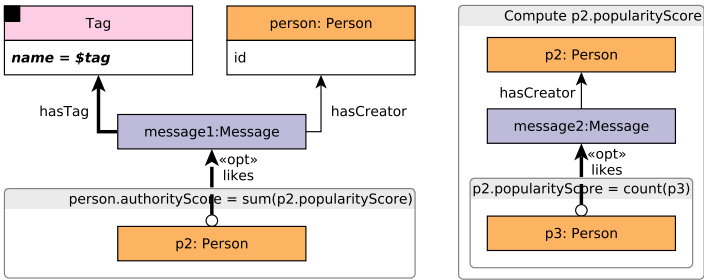
BI / read / 5

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

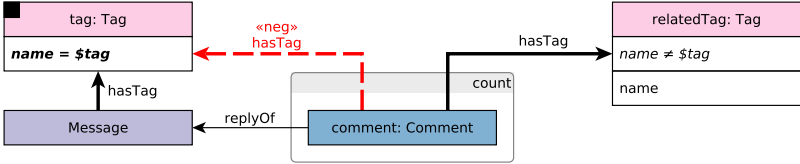
query	BI / read / 5																												
title	Most active posters of a given topic																												
pattern																													
desc.	<p>Get each Person (person) who has created a Message (message) with a given Tag (direct relation, not transitive). Considering only these Messages, for each Person node:</p> <ul style="list-style-type: none">Count its Messages (messageCount).Count likes (likeCount) to its Messages.Count Comments (replyCount) in reply to it Messages. <p>The score is calculated according to the following formula: $1 \times \text{messageCount} + 2 \times \text{replyCount} + 10 \times \text{likeCount}$.</p>																												
params	<div><div>1</div><div>tag</div><div>Long String</div></div>	Tags with a similar amount of Messages are selected. To avoid caching, different Tags should be used than the ones in Q6 and Q7.																											
result	<table><tr><td>1</td><td>person.id</td><td>ID</td><td>R</td><td></td></tr><tr><td>2</td><td>replyCount</td><td>32-bit Integer</td><td>A</td><td></td></tr><tr><td>3</td><td>likeCount</td><td>32-bit Integer</td><td>A</td><td></td></tr><tr><td>4</td><td>messageCount</td><td>32-bit Integer</td><td>A</td><td></td></tr><tr><td>5</td><td>score</td><td>32-bit Integer</td><td>A</td><td></td></tr></table>				1	person.id	ID	R		2	replyCount	32-bit Integer	A		3	likeCount	32-bit Integer	A		4	messageCount	32-bit Integer	A		5	score	32-bit Integer	A	
1	person.id	ID	R																										
2	replyCount	32-bit Integer	A																										
3	likeCount	32-bit Integer	A																										
4	messageCount	32-bit Integer	A																										
5	score	32-bit Integer	A																										
sort	<table><tr><td>1</td><td>score</td><td>↓</td><td></td></tr><tr><td>2</td><td>person.id</td><td>↑</td><td></td></tr></table>				1	score	↓		2	person.id	↑																		
1	score	↓																											
2	person.id	↑																											
limit	100																												
CPs	1.2, 2.3, 8.2																												

BI / read / 6

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 6			
title	Most authoritative users on a given topic			
pattern				
desc.	<p>Given a Tag (tag), find all Persons (person) that ever created a Message with the Tag. For each of these Persons (person) compute their “authority score” as follows:</p> <ul style="list-style-type: none">• The “authority score” is the sum of “popularity scores” of the Persons (p2) that liked any of that Person’s Messages with the given Tag (same criterion as for message1).• A Person’s (p2) “popularity score” is defined as the total number of likes on all of their Messages (message2).			
params	1	tag	Long String	Tags with a similar amount of Messages are selected. To avoid caching, different Tags should be used than the ones in Q5 and Q7.
result	1	person.id	ID	R
	2	authorityScore	32-bit Integer	A
sort	1	authorityScore	↓	
	2	person1.id	↑	
limit	100			
CPs	1.2, 2.3, 3.3, 6.1, 8.2			
relevance	Computing the authority scores might involve computing the popularity score for the same Person multiple times. Implementations are advised to avoid such redundant computations.			

BI / read / 7

BI 1	query	BI / read / 7			
BI 2	title	Related topics			
BI 3	pattern				
BI 4					
BI 5					
BI 6					
BI 7					
BI 8	desc.	Find all Messages that have a given Tag. Find the related Tags attached to (direct) reply Comments of these Messages, but only of those reply Comments that do not have the given Tag. Group the Tags by name, and get the count of replies in each group.			
BI 9					
BI 10					
BI 11					
BI 12					
BI 13	params	1	tag	Long String	Tags with a similar amount of Messages are selected. To avoid caching, different Tags should be used than the ones in Q5 and Q6.
BI 14					
BI 15					
BI 16					
BI 17					
BI 18	result	1	relatedTag.name	Long String	R
BI 19		2	count	32-bit Integer	A
BI 20					
	sort	1	count	↓	
		2	relatedTag.name	↑	
	limit	100			
	CPs	1.4, 3.3, 5.2, 8.1			

BI / read / 8

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 8			
title	Central person for a tag			
pattern	<div><div>For each person with a matching hasInterest and/or hasCreator edge, compute person.score = (if hasInterest edge exists then 100 else 0) + count(message)</div><div><pre>graph LR Tag[Tag] -- hasTag --> Message[message: Message] Person1[person: Person] -- hasInterest --> Tag Person1 -- hasCreator --> Message Person1 -- knows --> Friend[friend: Person] subgraph Count Message end subgraph FriendsScoreCalc Person1 -- "sum(friend.score)" --> FriendsScore[] end</pre></div></div>			
desc.	<p>Given a Tag, find all Persons that are interested in the Tag and/or have written a Message (Post or Comment) with a creationDate after a given date and that has a given Tag. For each Person, compute the score as the sum of the following two aspects:</p> <ul style="list-style-type: none">• 100, if the Person has this Tag as their interest, or 0 otherwise• number of Messages by this Person with the given Tag <p>Also, for each Person, compute the sum of the score of the Person’s friends (friendsScore).</p>			
params	<div><div>1</div><div>tag</div><div>Long String</div></div>	<div>Tags with a similar amount of Messages are selected</div> <div>(A): A range during which a flashmob event happened (it should yield at least a 5× difference)</div> <div>(B): A regular range (does not include a flashmob event)</div>		
result	<div><div>1</div><div>person.id</div><div>ID</div><div>R</div></div> <div><div>2</div><div>score</div><div>32-bit Integer</div><div>A</div></div> <div><div>3</div><div>friendsScore</div><div>32-bit Integer</div><div>A</div></div>	<div>The sum of the score of the person’s friends</div>		
sort	<div><div>1</div><div>score + friendsScore</div><div>↓</div></div> <div><div>2</div><div>person.id</div><div>↑</div></div>			
limit	100			
CPs	1.2, 2.1, 2.3, 3.2, 5.3, 8.2, 8.4, 8.5			
relevance	Similarly to BI 16, there are two major ways to compute this query: (1) creating an induced subgraph of the interested Persons and their friends and performing the scoring on this graph or (2) performing the scoring without creating an induced subgraph and scoring the friends of a Person on-the-fly. The first approach is more efficient as it avoids redundant computations, however, specifying it needs support for composable graph queries.			

BI / read / 9

BI 1

BI 2

BI 3

BI 4

BI 5

BI 6

BI 7

BI 8

BI 9

BI 10

BI 11

BI 12

BI 13

BI 14

BI 15

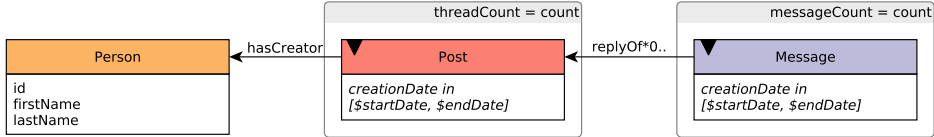
BI 16

BI 17

BI 18

BI 19

BI 20

query	BI / read / 9				
title	Top thread initiators				
pattern					
desc.	<p>For each Person, count the number of Posts they created in the time interval [startDate, endDate] (equivalent to the number of threads they initiated) and the number of Messages in each of their (transitive) reply trees, including the root Post of each tree. When calculating Message counts only consider Messages created within the given time interval.</p> <p>Return each Person, number of Posts they created, and the count of all Messages that appeared in the reply trees (including the Post at the root of tree).</p>				
params	1	startDate	Date	Selected around the same date	
	2	endDate	Date	80-100 days after the startDate	
result	1	person.id	ID	R	
	2	person.firstName	String	R	
	3	person.lastName	String	R	
	4	threadCount	32-bit Integer	A	The number of Posts created by that Person (the number of threads initiated)
	5	messageCount	32-bit Integer	A	The number of Messages created in all the threads this Person initiated
sort	1	messageCount	↓		
	2	person.id	↑		
limit	100				
CPs	1.2, 2.2, 2.3, 3.2, 7.2, 7.3, 7.4, 8.1, 8.5				

BI / read / 10

BI 1	query	BI / read / 10				
BI 2	title	Experts in social circle				
BI 3	pattern					
BI 4						
BI 5						
BI 6						
BI 7						
BI 8						
BI 9						
BI 10						
BI 11						
BI 12						
BI 13	desc.	<p>Given a Person (startPerson), find all other Persons (expertCandidatePerson) that live in a given Country and are connected to given Person by a <i>shortest path</i> with length in range [minPathDistance, maxPathDistance] through the knows relation.</p> <p>For each of these expertCandidatePerson nodes, retrieve all of their Messages that contain at least one Tag belonging to a given TagClass (direct relation not transitive). For each Message, retrieve all of its Tags.</p> <p>Group the results by Persons and Tags, then count the Messages by a certain Person having a certain Tag.</p>				
BI 16						
BI 17						
BI 18						
BI 19	params	1	personId	ID	(A) Persons with an average degree of knows edges are selected (B) Persons who have only one friend and that Person has two friends in total (including the original Person)	
BI 20		2	country	String	Select mid-sized Countries	
		3	tagClass	Long String	TagClasses with a similar degree of hasType edges are selected	
		4	minPathDistance	32-bit Integer	3	
		5	maxPathDistance	32-bit Integer	4	
	result	1	expertCandidatePerson.id	ID	R	
		2	tag.name	Long String	R	
		3	messageCount	32-bit Integer	A	Number of Messages created by that Person containing that Tag
	sort	1	messageCount	↓		
		2	tag.name	↑		
		3	expertCandidatePerson.id	↑		
	limit	100				
	CPs	1.2, 1.3, 2.3, 2.4, 3.3, 5.3, 7.1, 7.2, 7.3, 8.1, 8.6				

BI / read / 11

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 11			
title	Friend triangles			
pattern				
desc.	<p>For a given country, count all the distinct triples of Persons such that:</p> <ul style="list-style-type: none">• a is friend of b,• b is friend of c,• c is friend of a, <p>and these friendships were created in the range [startDate, endDate].</p> <p>Distinct means that given a triple t_1 in the result set R of all qualified triples, there is no triple t_2 in R such that t_1 and t_2 have the same set of elements.</p>			
params	1	country	Long String	Selected from the largest Countries (India, China)
	2	startDate	Date	Selected from a 30-day interval towards the end of the simulation time
	3	endDate	Date	Selected to yield around a 100-day interval
result	1	count	64-bit Integer	A
limit	n/a			
CPs	1.1, 2.3, 2.5			

BI / read / 12

BI 1	query	BI / read / 12			
BI 2	title	How many persons have a given number of messages			
BI 3	pattern				
BI 4					
BI 5					
BI 6					
BI 7					
BI 8					
BI 9					
BI 10	desc.	For each Person, count the number of Messages they made (messageCount). Only count Messages with the following attributes:			
BI 11		<ul style="list-style-type: none">• Its content is not empty (and consequently, the imageFile attribute is empty for Posts).• Its length is below the lengthThreshold (exclusive, equality is not allowed).• Its creationDate is after startDate (exclusive, equality is not allowed).• It is written in any of the given languages.			
BI 12		<ul style="list-style-type: none">– The language of a Post is defined by its language attribute.– The language of a Comment is that of the Post that initiates the thread where the Comment replies to.			
BI 13		The Post and Comments in the reply tree’s path (from the Message to the Post) do not have to satisfy the constraints for content, length, and creationDate.			
BI 14		For each messageCount value, count the number of Persons with exactly messageCount Messages (with the required attributes).			
BI 15					
BI 16					
BI 17					
BI 18					
BI 19					
BI 20					
params	1	startDate	Date	Selected randomly from a 60-day interval.	
	2	lengthThreshold	32-bit Integer	Balanced against startDate to filter around 30% of the Messages within a language and keep the variance low. The selection of this parameter uses a factor table of bucketed Message lengths and creation dates.	
	3	languages	{String}	Only the most frequently used languages	
result	1	messageCount	32-bit Integer	A	Number of Messages created
	2	personCount	32-bit Integer	A	Number of Persons with messageCount Messages
sort	1	personCount	↓		
	2	messageCount	↓		
limit	n/a				
CPs	1.1, 1.2, 1.4, 3.2, 4.2, 4.3, 8.1, 8.2, 8.3, 8.4, 8.5				

BI / read / 13

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 13																								
title	Zombies in a country																								
pattern	<div><div>1. zombies = collect(zombie)</div><div><div><div>Country</div><div><div>name = \$country</div></div></div><div><div>City</div><div><div>isPartOf</div><div>Country</div></div><div><div>isLocatedIn</div><div>zombie: Person</div></div></div><div><div>zombie: Person</div><div><div>creationDate < \$endDate and (messageCount / months < 1)</div></div></div><div><div>message: Message</div><div><div>messageCount = count(message)</div><div>creationDate < \$endDate</div></div></div><div><div>zombie: Person</div><div><div>«opt» hasCreator</div><div>message: Message</div></div></div></div><div><div>2. For each zombie IN zombies, calculate: zombieScore = zombieLikeCount / totalLikeCount</div><div><div><div>likerPerson: Person</div><div><div>totalLikeCount = count(likerPerson)</div><div>creationDate < \$endDate</div></div></div><div><div>zombie: Person</div><div><div>hasCreator</div><div>Message</div></div></div><div><div>likerZombie: Person</div><div><div>zombieLikeCount = count(likerZombie)</div><div>creationDate < \$endDate and likerZombie IN zombies</div></div></div><div><div>likerPerson: Person</div><div><div>«opt» likes</div><div>Message</div></div></div><div><div>Message</div><div><div>«opt» likes</div><div>likerZombie: Person</div></div></div></div></div></div>																								
desc.	<p>Find zombies within the given country, and return their zombie scores. A zombie is a Person created before the given endDate, which has created an average of [0, 1) Messages per month, during the time range between profile’s creationDate and the given endDate. The number of months spans the time range from the creationDate of the profile to the endDate with partial months on both end counting as one month (e.g. a creationDate of Jan 31 and an endDate of Mar 1 result in 3 months).</p> <p>For each zombie, calculate the following:</p> <ul style="list-style-type: none">zombieLikeCount: the number of likes received from other zombies.totalLikeCount: the total number of likes received.zombieScore: zombieLikeCount / totalLikeCount. If the value of totalLikeCount is 0, the zombieScore of the zombie should be 0.0. <p>For both zombieLikeCount and totalLikeCount, only consider likes received from profiles that were created before the given endDate.</p>																								
params	<table><tr><td>1</td><td>country</td><td>Long String</td><td colspan="2">Selected from the largest Countries (India, China)</td></tr><tr><td>2</td><td>endDate</td><td>Date</td><td colspan="2">Selected from the last days of the initial data set</td></tr></table>					1	country	Long String	Selected from the largest Countries (India, China)		2	endDate	Date	Selected from the last days of the initial data set											
1	country	Long String	Selected from the largest Countries (India, China)																						
2	endDate	Date	Selected from the last days of the initial data set																						
result	<table><tr><td>1</td><td>zombie.id</td><td>ID</td><td>R</td><td></td></tr><tr><td>2</td><td>zombieLikeCount</td><td>32-bit Integer</td><td>A</td><td></td></tr><tr><td>3</td><td>totalLikeCount</td><td>32-bit Integer</td><td>A</td><td></td></tr><tr><td>4</td><td>zombieScore</td><td>32-bit Float</td><td>A</td><td>Determined as zombieLikeCount / totalLikeCount</td></tr></table>					1	zombie.id	ID	R		2	zombieLikeCount	32-bit Integer	A		3	totalLikeCount	32-bit Integer	A		4	zombieScore	32-bit Float	A	Determined as zombieLikeCount / totalLikeCount
1	zombie.id	ID	R																						
2	zombieLikeCount	32-bit Integer	A																						
3	totalLikeCount	32-bit Integer	A																						
4	zombieScore	32-bit Float	A	Determined as zombieLikeCount / totalLikeCount																					
sort	<table><tr><td>1</td><td>zombieScore</td><td>↓</td><td></td></tr><tr><td>2</td><td>zombie.id</td><td>↑</td><td></td></tr></table>					1	zombieScore	↓		2	zombie.id	↑													
1	zombieScore	↓																							
2	zombie.id	↑																							
limit	100																								
CPs	1.2, 2.1, 2.3, 2.4, 3.2, 3.3, 4.2, 5.1, 5.3, 8.2, 8.4, 8.5																								

BI / read / 14

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 14			
title	International dialog			
pattern	<div><div>For each pair of countries, calculate the cost as a sum of cases #1-4. Cases that have a match add to the final score with the specified value. Each case only counts once, multiple matches do not increase to the score.</div><div><div><div><div>Country</div><div><div>name = \$country1</div></div></div><div><div>city1: City</div><div><div>name</div></div></div><div><div>person1: Person</div><div><div>id</div></div></div><div><div>isPartOf</div></div><div><div>isLocatedIn</div></div></div><div><div><div>Country</div><div><div>name = \$country2</div></div></div><div><div>City</div><div></div></div><div><div>person2: Person</div><div><div>id</div></div></div><div><div>isPartOf</div></div><div><div>isLocatedIn</div></div><div><div>knows</div></div></div></div><div><div>Case 1: score += 4</div><div><div><div>person1: Person</div><div>hasCreator</div><div>Comment</div></div><div><div>person2: Person</div><div>hasCreator</div><div>Message</div></div><div><div>replyOf</div></div></div></div><div><div>Case 2: score += 1</div><div><div><div>person1: Person</div><div>hasCreator</div><div>Message</div></div><div><div>person2: Person</div><div>hasCreator</div><div>Comment</div></div><div><div>replyOf</div></div></div></div><div><div>Case 3: score += 10</div><div><div><div>person1: Person</div><div>likes</div><div>Message</div></div><div><div>person2: Person</div><div>hasCreator</div><div>Message</div></div></div></div><div><div>Case 4: score += 1</div><div><div><div>person1: Person</div><div>hasCreator</div><div>Message</div></div><div><div>person2: Person</div><div>likes</div><div>Message</div></div></div></div></div>			
desc.	<p>Consider all pairs of people (person1, person2) such that (1) they know each ther, (2) one is located in a City of Country country1, and (3) the other is located in a City of Country country2. For each City of Country country1, return the highest scoring pair. The score of a pair is defined as the sum of the subscores awarded for the following kinds of interaction. The initial value is score = 0.</p> <ol style="list-style-type: none">person1 has created a reply Comment to at least one Message by person2: score += 4person1 has created at least one Message that person2 has created a reply to: score += 1person1 liked at least one Message by person2: score += 10person1 has created at least one Message that was liked by person2: score += 1 <p>Consequently, the maximum score a pair can obtain is: 4 + 1 + 10 + 1 = 16.</p>			
params	<div><div><div>1</div><div>country1</div><div>Long String</div></div><div><div>2</div><div>country2</div><div>Long String</div></div></div>	<div>(A) Correlated with parameter country2, i.e. the Countries are close and there are many Persons knowing each other</div> <div>(B) Uncorrelated with parameter country2, i.e. the Countries are afar and there are few Persons knowing each other</div>		
result	<div><div><div>1</div><div>person1.id</div><div>ID</div><div>R</div></div><div><div>2</div><div>person2.id</div><div>ID</div><div>R</div></div><div><div>3</div><div>city1.name</div><div>Long String</div><div>R</div></div><div><div>4</div><div>score</div><div>32-bit Integer</div><div>C</div></div></div>			
sort	<div><div><div>1</div><div>score</div><div>↓</div></div><div><div>2</div><div>person1.id</div><div>↑</div></div><div><div>3</div><div>person2.id</div><div>↑</div></div></div>			
limit	n/a			
CPs	1.3, 1.4, 2.1, 3.1, 3.3, 5.1, 5.2, 5.3, 8.3, 8.4			

BI / read / 15

BI 1	query	BI / read / 15			
BI 2	title	Trusted connection paths through forums created in a given timeframe			
BI 3	pattern	<div><div>Enumerate all unweighted shortest paths on knows edges between person1 to person2.</div><div><div>person1: Person</div><div>knows*</div><div>person2: Person</div><div><div>id = \$person1Id</div><div>id = \$person2Id</div></div></div></div>			
BI 4		<div><div>For each knows edge in the path, calculate a weight based on interactions between the pair of Persons of the edge, calculated as a sum of cases #1 and #2 for the Persons (both ways), and the sum of these weights determine the total weight of each path.</div><div><div>p1</div><div>knows</div><div>pX</div><div>knows</div><div>pY</div><div>...</div><div>pW</div><div>knows</div><div>p2</div></div></div>			
BI 5		<div><div>Case 1: Replies on Posts, weight += 1.0 × count(c)</div><div><div>personA: Person</div><div>knows</div><div>personB: Person</div><div><div>hasCreator</div><div>c: Comment</div><div>replyOf</div><div>post: Post</div><div><div>containerOf</div><div>forum: Forum</div><div>creationDate in [\$startDate, \$endDate]</div></div></div></div></div>			
BI 6		<div><div>Case 2: Replies on Comments, weight += 0.5 × count(c1)</div><div><div>personA: Person</div><div>knows</div><div>personB: Person</div><div><div>hasCreator</div><div>c1: Comment</div><div>replyOf</div><div>c2: Comment</div><div><div>replyOf*</div><div>Post</div><div><div>containerOf</div><div>forum: Forum</div><div>creationDate in [\$startDate, \$endDate]</div></div></div></div></div></div>			
BI 7					
BI 8					
BI 9					
BI 10					
BI 11					
BI 12					
BI 13					
BI 14					
BI 15					
BI 16					
BI 17					
BI 18		desc.	<p>Given two Persons, find all (unweighted) shortest paths between these two Persons, in the subgraph induced by the knows relationship.</p> <p>Then, for each path calculate a weight. The nodes in the path are Persons, and the weight of a path is the sum of weights between every pair of consecutive Person nodes in the path.</p> <p>The weight for a pair of Persons is calculated based on their interactions:</p> <ul style="list-style-type: none">• Every direct reply (by one of the Persons) to a Post (by the other Person) contributes 1.0.• Every direct reply (by one of the Persons) to a Comment (by the other Person) contributes 0.5. <p>Only consider Messages that were created in a Forum that was created within the timeframe (interval) [startDate, endDate]. Note that for Comments, the containing Forum is that of the Post that the comment (transitively) replies to. Also note that interactions are counted both ways.</p> <p>Return all paths with the Person IDs ordered by their weights descending.</p>		
BI 19					
BI 20					
params	1	person1Id	ID	(A) person1Id – person2Id pair with a distance of exactly 4 hops (B) person1Id – person2Id pair with a distance of exactly 2 hops	
	2	person2Id	ID		
	3	startDate	Date	(A) Small interval (approx. one week) (B) Big interval (approx. one month)	
	4	endDate	Date		
result	1	person.id	[ID]	C	Ordered sequence of the Person IDs in the path
	2	weight	32-bit Float	C	
sort	1	weight	↓	The order of paths with the same weight is unspecified	
	2	personIds	↑	The IDs in the paths are used for lexicographical sorting	
limit	n/a				
CPs	1.2, 2.1, 2.2, 2.4, 3.3, 5.1, 5.3, 7.2, 7.3, 7.5, 7.7, 8.1, 8.2, 8.3, 8.4, 8.5, 8.6				

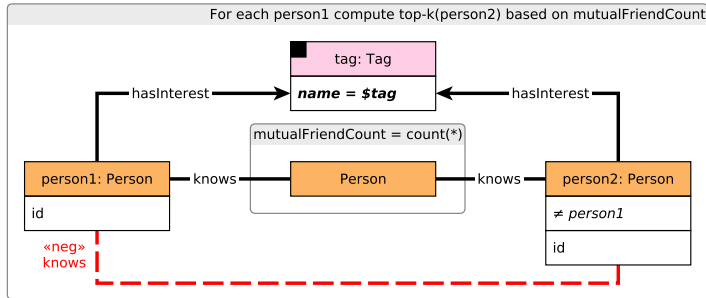
BI / read / 16

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 16			
title	Fake news detection			
pattern	<div><div>For \$tagX/\$dayX in [tagA/dateA, tagB/dateB], compute scoreX = count(messageX)</div><div><div>1. Create an induced subgraph of Persons who created a Message with Tag \$tagX on \$dateX</div><div><div><div>tag: Tag</div><div>name = \$tagX</div></div><div>hasTag</div><div><div>Message</div><div>day(creationDate) = \$dateX</div></div><div>hasCreator</div><div><div>person: Person</div></div></div></div><div><div>2. In the subgraph, count the Messages (using the same conditions) from People with ≤ \$maxKnowsLimit friends</div><div><div><div>tag: Tag</div><div>name = \$tagX</div></div><div>hasTag</div><div><div>messageX: Message</div><div>day(creationDate) = \$dateX</div></div><div>hasCreator</div><div><div>person: Person</div><div>count ≤ \$maxKnowsLimit</div><div>«opt» knows</div><div>Person</div></div></div></div></div>			
desc.	<p>Given two Tag/date pairs (tagA/dateA and tagB/dateB), for each pair tagX/dateX:</p> <ul style="list-style-type: none">• Create an induced subgraph between Persons where for each pair of Persons person1/person2, both have created a Message on the day of dateX with Tag tagX.• In the induced subgraph, only keep pairs of Persons who have at most maxKnowsLimit friends (in the induced subgraph).• For these Persons, count the number of Messages created on dateX with Tag tagX. <p>Return Persons who had at least one Messages for both tagA/dateA and tagB/dateB ranked by their total number of Messages (descending).</p>			
params	<div><div>1</div><div>tagA</div><div>Long String</div></div> <div><div>2</div><div>dateA</div><div>Date</div></div> <div><div>3</div><div>tagB</div><div>Long String</div></div> <div><div>4</div><div>dateB</div><div>Date</div></div> <div><div>5</div><div>maxKnowsLimit</div><div>32-bit Integer</div></div>	<div>(A) tagA–dateA, tagB–dateB are both selected to be a flashmob Tag – date combination</div> <div>(B) tagA–dateA, tagB–dateB are both selected to be a non-flashmob Tag – date combination</div>		
result	<div><div>1</div><div>person.id</div><div>ID</div></div> <div><div>2</div><div>messageCountA</div><div>32-bit Integer</div></div> <div><div>3</div><div>messageCountB</div><div>32-bit Integer</div></div>	<div>R</div> <div>A</div> <div>A</div>	<div></div> <div>Message count for tagA/dateA</div> <div>Message count for tagB/dateB</div>	
sort	<div><div>1</div><div>messageCountA + messageCountB</div><div>↓</div></div> <div><div>2</div><div>person.id</div><div>↑</div></div>			
limit	20			
CPs	5.3, 8.4, 8.5			
relevance	There are two major ways to compute this query: (1) create the induced subgraph as suggested by the specification (either as a view or in materialized form), or (2) skip creating the induced subgraph and perform on-the-fly check for the number of friends (who also posted at least one Message with the given Tag on the given date). The latter approach is easier to express in systems which do not provide graph views but might result in redundant computations (the query engine might repeatedly check whether a Person has at least one Message that satisfies the conditions).			

BI / read / 18

BI 1
BI 2
BI 3
BI 4
BI 5
BI 6
BI 7
BI 8
BI 9
BI 10
BI 11
BI 12
BI 13
BI 14
BI 15
BI 16
BI 17
BI 18
BI 19
BI 20

query	BI / read / 18				
title	Friend recommendation				
pattern	<div>For each person1 compute top-k(person2) based on mutualFriendCount</div> 				
desc.	<div>For a given Tag (tag), for each person1 interested in tag, recommend new friends (person2) who</div> <ul style="list-style-type: none">• do not yet know person1• at least one mutual friend with person1• are also interested in tag. <div>Rank Persons person2 based on the number of mutual friends with person1.</div>				
params	1	tag	Long String	Tags with a similar amount of Persons are selected	
result	1	person1.id	ID	R	
	2	person2.id	ID	R	
	3	mutualFriendCount	32-bit Integer	A	
sort	1	mutualFriendCount	↓		
	2	person1.id	↑		
	3	person2.id	↑		
limit	20				
CPs	2.5, 8.1				

BI / read / 19

BI 1	query	BI / read / 19			
BI 2	title	Interaction path between cities			
BI 3	pattern	<p>Find the shortest paths between all pairs of Persons in city1 and city2</p> <p>city1: City $id = \\$city1id$ isLocatedIn person1: Person</p> <p>city2: City $id = \\$city2id$ isLocatedIn person2: Person</p> <p>compute weighted shortest paths on knows.weight</p> <p>The weight of a knows edge is based on the number of interactions between its Persons: $knows.weight = 1 / (count(i1) + count(i2))$</p> <p>Case i1: Reply from personA to Person B's Message</p> <p>Case i2: Reply from personB to personA's Message</p>			
BI 13	desc.	<p>Given two Cities city1, city2, find Persons person1, person2 living in these Cities (respectively) with the shortest <i>interaction path</i> between them. If there are multiple pairs of people with shortest paths having the same total weight, return all of them.</p> <p>The shortest path is computed using a weight between two Persons defined as the reciprocal of the number of interactions (direct reply Comments to a Message by the other Person). Therefore, more interactions imply a smaller weight.</p> <p><i>Note:</i> Interactions are counted both ways, i.e. if Alice writes 2 reply Comments to Bob's Messages and Bob writes 3 reply Comments to Alice's Messages, their total number of interactions is 5.</p>			
BI 14	params	1	city1Id	ID	(A) Small Cities within the same Country with many direct relationships between their inhabitants
BI 15		2	city2Id	ID	(B) Small Cities from different Countries with only a few direct relationships between their inhabitants
BI 16	result	1	person1.id	ID	R
BI 17		2	person2.id	ID	R
BI 18		3	totalWeight	32-bit Float	C
BI 19	sort	1	totalWeight	↑	
BI 20		2	person1.id	↑	
		3	person2.id	↑	
	limit	20			
	CPs	3.3, 7.6, 7.7, 8.4, 8.6			
	relevance	<p>Finding shortest paths between pairs of Persons in Cities can be implemented in theory with an <i>all-pairs shortest paths</i> algorithm. However, this needs to be executed on the whole Person-knows-Person graph (with edge weights derived from the number of interactions) so it is expected to be prohibitively expensive. A better approach is using multiple <i>single-source shortest path algorithms</i> (e.g. from the City with fewer inhabitants). Implementations can either pre-compute edge weights or compute them on-the-fly.</p>			

BI / read / 20

BI 1	query	BI / read / 20			
BI 2	title	Recruitment			
BI 3	pattern				
BI 4					
BI 5					
BI 6					
BI 7					
BI 8	desc.	Given a Company company and a Person person2 (who is not working and has not worked at company), find a different Person (person1) who works or at some point worked in company and is reachable by from person2 through people who have studied together. On this path, we only consider edges between Persons who know each other and attended the same University and set the weight of the edge to the absolute difference between the year of enrolment plus 1 (studyAt.classYear + 1). If the Persons attended multiple universities, we select the smallest (min) value. If there are multiple Person person1 nodes with the same shortest path, return all of them.			
BI 9					
BI 10					
BI 11					
BI 12					
BI 13	params	1	company	Long String	Companies with a similar number of employees (former or current) are selected
BI 14		2	person2Id	ID	person2 is selected so that there is no direct (1-hop) path to any person1 working at company
BI 15					
BI 16					
BI 17					
BI 18	result	1	person1.id	ID	R
BI 19		2	totalWeight	64-bit Integer	C
BI 20	sort	1	totalWeight	↑	
		2	person1.id	↑	
	limit	20			
	CPs	3.3, 7.6, 7.7, 8.4, 8.6			
	relevance	Implementations can either pre-compute edge weights or compute them on-the-fly. To find the (weighted) shortest path efficiently, can use e.g. a bidirectional Dijkstra algorithm.			

1.5 Refreshes

1.5.1 Inserts

See *Interactive Inserts* (??).

1.5.2 Deletes


Each delete query removes

1. a single edge between two existing nodes
2. or a node, all incident edges and, in certain cases, nodes and edges that are transitively reachable on a certain path.


BI / delete / 1

DEL 1	query	BI / delete / 1				
DEL 2	title	Remove person and its personal forums and message (sub)threads				
DEL 3	pattern					
DEL 4						
DEL 5						
DEL 6						
DEL 7						
DEL 8						
	desc.	Remove a Person and its edges (isLocatedIn, studyAt, workAt, hasInterest, likes, knows, hasMember, hasModerator, hasCreator). Additionally, remove the Album and Wall Forums whose moderator is the Person and remove all Messages the Person has created in the rest of the Forums (Groups).				
	params	<table><tr><td>1</td><td>personId</td><td>ID</td><td></td></tr></table>	1	personId	ID	
1	personId	ID				
	CPs	9.3, 9.4, 9.5				
	relevance	<ul style="list-style-type: none">Removal of a Person removes Forums of type “Walls” and “Albums” but not “Groups”, which can continue if even the founder has left the network. For Groups, the hasModerator edge is deleted. We have discussed various approaches to appoint a new moderator, e.g.<ol style="list-style-type: none">choose member at random from the set of existing group members orthe member with the oldest group join date becomes the moderator. However, to keep the generator and the workload simple, currently no moderator is selected, leaving the group without a moderator.Removal of a Person removes all Posts/Comments they are creator of this could result in the removal of a Comment in the middle of a thread.				

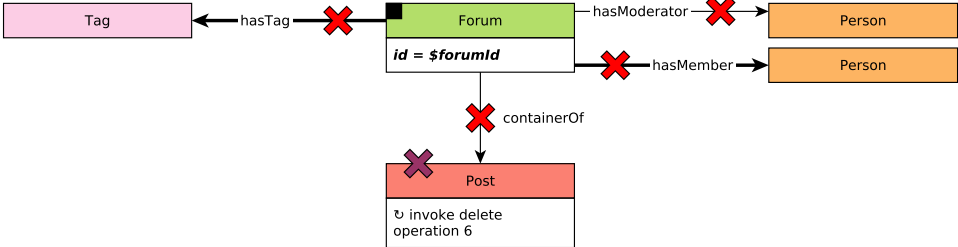
BI / delete / 2

DEL 1	query	BI / delete / 2			
DEL 2	title	Remove post like			
DEL 3	pattern				
DEL 4					
DEL 5					
DEL 6					
DEL 7	desc.	Given a Person and a Post, remove the likes edge between them.			
DEL 8	params	<div>1</div>	personId	ID	
		<div>2</div>	postId	ID	
	CPs	9.4			
	relevance	Removal of a likes edge is a rare event, e.g. people accidentally liking a Post, this can be reflected by the relative frequency of the operation.			

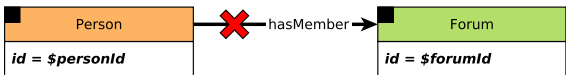
BI / delete / 3

DEL 1	query	BI / delete / 3		
DEL 2	title	Remove comment like		
DEL 3	pattern			
DEL 4				
DEL 5				
DEL 6				
DEL 7	desc.	Given a Person and a Comment, remove the likes edge between them.		
DEL 8	params	1	personId	ID
		2	commentId	ID
	CPs	9.4		
	relevance	Removal of a likes edge is a rare event, e.g. people accidentally liking a Comment, this can be reflected by the relative frequency of the operation.		

BI / delete / 4

DEL 1	query	BI / delete / 4		
DEL 2	title	Remove forum and its content		
DEL 3	pattern			
DEL 4				
DEL 5				
DEL 6				
DEL 7	desc.	Remove a Forum and its edges (hasModerator, hasMember, hasTag) and all Posts in the Forum (connected by containerOf edges) and their direct and transitive Comments.		
DEL 8	params	1	forumId	ID
	CPs	9.3, 9.4, 9.5		
	relevance	n/a		

BI / delete / 5

DEL 1	query	BI / delete / 5		
DEL 2	title	Remove forum membership		
DEL 3	pattern			
DEL 4				
DEL 5				
DEL 6				
DEL 7	desc.	Given a Forum and a Person, remove the hasMember edge between them.		
DEL 8	params	1	forumId	ID
		2	personId	ID
	CPs	9.4		
	relevance	n/a		

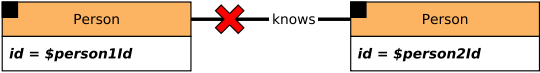
BI / delete / 6

DEL 1	query	BI / delete / 6	
DEL 2	title	Remove post thread	
DEL 3	pattern		
DEL 4			
DEL 5			
DEL 6			
DEL 7			
DEL 8			
	desc.	Remove a Post node and its edges (isLocatedIn, likes, hasCreator, hasTag, containerOf). Remove all replies to the Post and the connecting replyOf edges. In addition, remove all transitive reply Comments to the Post and their edges.	
	params	1	postId ID
	CPs	9.3, 9.4, 9.5	
	relevance	n/a	

BI / delete / 7

DEL 1	query	BI / delete / 7	
DEL 2	title	Remove comment subthread	
DEL 3	pattern		
DEL 4			
DEL 5			
DEL 6			
DEL 7			
DEL 8			
	desc.	Remove a Comment node and its edges (isLocatedIn, likes, hasCreator, hasTag). In addition, remove all replies to the Comment connected by replyOf and their edges.	
	params	1	commentId ID
	CPs	9.3, 9.4, 9.5	
	relevance	n/a	

BI / delete / 8

DEL 1	query	BI / delete / 8		
DEL 2	title	Remove friendship		
DEL 3	pattern			
DEL 4				
DEL 5				
DEL 6				
DEL 7				
DEL 8	desc.	Given two Person nodes, remove the knows edge between them.		
	params	1	person1Id	ID
		2	person2Id	ID
	CPs	9.4		
	relevance	n/a		

BIBLIOGRAPHY

- [1] Andrey Gubichev and Peter A. Boncz. “Parameter Curation for Benchmark Queries”. In: *TPCTC*. Vol. 8904. Lecture Notes in Computer Science. Springer, 2014, pp. 113–129.
- [2] Gábor Szárnyas et al. “An early look at the LDBC Social Network Benchmark’s Business Intelligence workload”. In: *GRADES-NDA at SIGMOD/PODS*. ACM, 2018, 9:1–9:11. doi: 10.1145/3210259.3210268.