

# TENTAMEN (EXAMINATION)

7

Tentamensdatum/*Examination date:* 1/11-19  
(åå-mm-dd/yy-mm-dd)

AID-nummer  
*AID number*

Ifyller av student

1	2	8	7		
---	---	---	---	--	--

*Completed by student*

Ifyller av vakt

1	2	8	7		
---	---	---	---	--	--

*Completed by supervisor*

Utbildningskod/*Education code:* TDDE 31 Modul/*Module:* TEN1

Kursnamn/*Course title:* Bio Data Analytics

Institution/*Department:* IDA

Jag intygar att varken mobil eller något annat otillåtet hjälpmedel finns tillgängligt under tentamen.  
*I confirm that no mobile or other non-permitted aids are available during the examination.*

Inlämnat: antal lösblad 8 tentamensformulär   
*Enclosed: number of sheets* *exam booklet*

Markera behandlade uppgifter med X/*Mark tasks attempted with an X*

X här/here	1	2	3	4	5	6	7	8	9	10	11	12	13	14	15
Erhållna poäng <i>Points obtained</i>	2	2.5	2.5	1.5	2.5	2.5	—	2	—	3					
X här/here	16	17	18	19	20	21	22	23	24	25	26	27	28	29	30
Erhållna poäng <i>Points obtained</i>															

Anvisningar/*Instructions*

1. Skriv AID-nummer, datum, kurskod och provkod på varje blad som lämnas in/  
*Write AID number, date, course code and exam code on every sheet that is handed in*
2. På varje papper får högst en uppgift lösas om inget annat anges/  
*Maximum one task per sheet unless otherwise instructed*
3. Skriv endast på papprets ena sida om inget annat anges/  
*Use only one side of each sheet unless otherwise instructed*
4. Numrera de papper som lämnas in/*Number every sheet that is handed in*
5. Använd inte röd penna/*Do not use a red pen/pencil*

Sen inlämning  
*Late hand in*

Klockslag \_\_\_\_\_  
*Time*

Orsak \_\_\_\_\_  
*Reason*

Σ Poäng/*Points:* 18.5 Betyg/*Grade:* 3

Examinator/*Examiner:* \_\_\_\_\_

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN 1

Blad nummer:  
Sheet number:

1

1

Volume: The amount of data e.g. Wikipedia

Veriaty: Different types of data, for example:

structured: dates, time

unstructured: mix of picture and documents.

Veracity: The trustworthiness of the data. Can you trust the information?

Velocity: The frequency of data that has to be processed, e.g. finance Market or Social Media.

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19	Blad nummer: Sheet number:
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN1	2

2 a) Vertical scalability is upgrading the existing node.

With example more ram, more memory, better cpus.

Horizontal scalability is adding more nodes into a distributed system.

1/1

b) read scalability means that the system can handle an increased number of read operations without losing performance.

Write scalability means that the system can handle an increased number of write operations without losing performance.

1/1

c) Data scalability is important within for example a flight system. If the system where not able to handle an increased amount of data, it could be devastating.

What kind of  
data is increasing?

0.5/1

2.5/1

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN1

Blad nummer: Sheet number:
3

- 3 a) the typically implemented queries uses keys to retrieve value. So we would have to loop through the keys until we found the one we where searching for. ✓ // /
- b) This would be done by using "Alice" as key add a value with IDs that has Alice as name. "Alice" → " , [alice-in-se]" // /
- c) In a document database, the information about each of these persons would be its own document. In a key database only the key has to be unique. But in a document database the name describing the value all has to be unique. within a document.
- only  
one  
difference
- not necessarily 05/11
- (25P)

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN1

Blad nummer:  
Sheet number:

4

4

## BASE

Basically Available - The system is available basically all the time.

Soft state - The system ~~is always changing even if no input has been made.~~

Eventually Consistent - This means that if no new inputs are made, then eventually all nodes with that data will have the latest update.

1.5 p

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN1

Blad nummer: Sheet number:
5

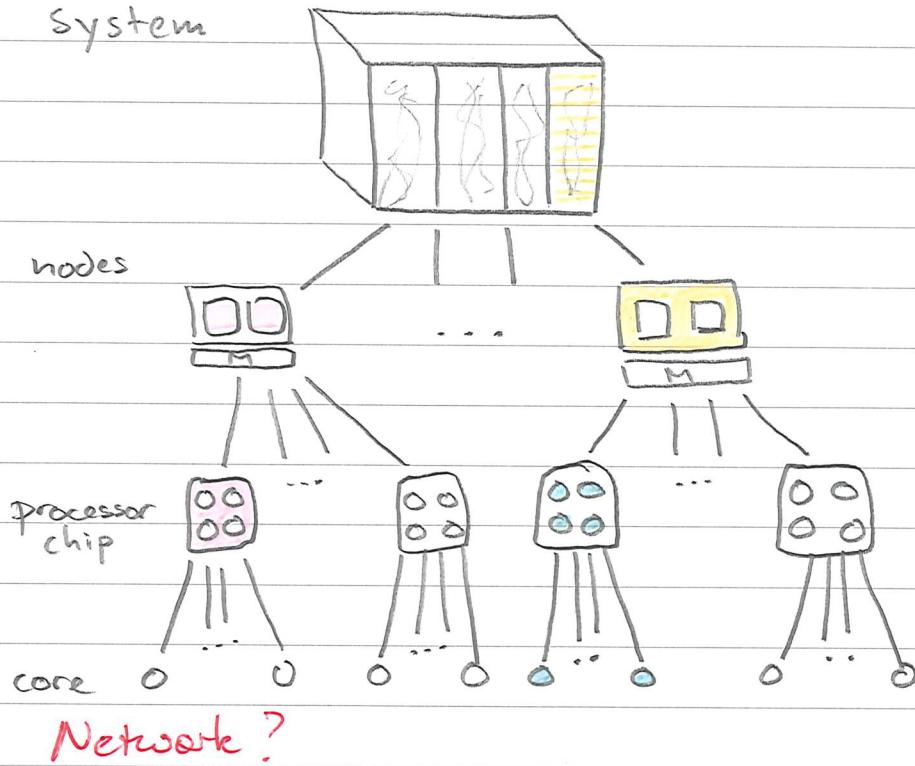
5

- a) a distributed file is divided into chunks stored separately. A namenode (Master) keeps track on where the parts of a file is saved and the information in each chunk will be handled by datanodes (Workers). Each chunk of a file will be saved at three different places to ensure Fault Tolerance. Multiple readers makes it good for parallelization. ✓

b) /

1,5

c) System



(✓)

0,75

Network?

- d) Makes the operations over data more efficient  
How?

0

2,25

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19	Blad nummer: Sheet number:
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN1	6

6

a) **Who?** should not depend on order between rows? of information since the order might differ. **where?** 0

b) first step is the read step, since the collection could be done at a different place this could be network I. ✓

(fourth step is the partitioning. Which) **Step 5. (shuffle)** sends the information to the node for reducing phase. which could involve network I/O ✓

Seventh step is the write which could also involve network O. ✓

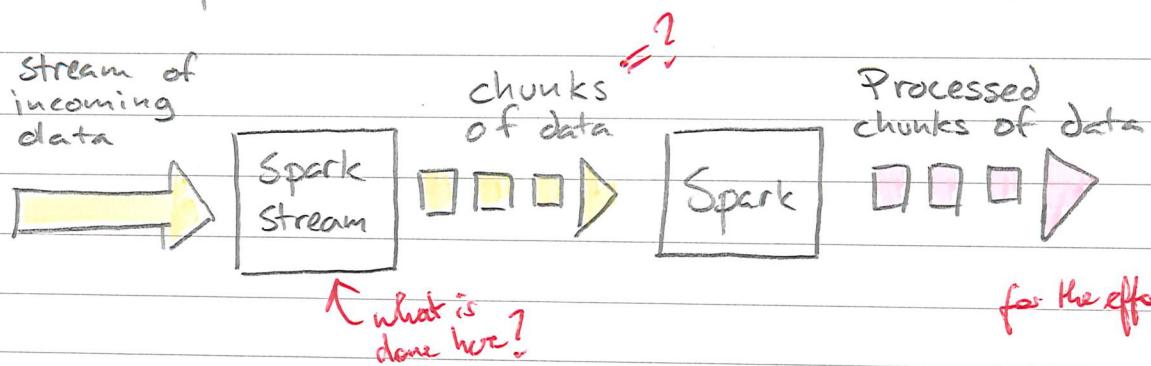
c) Could be beneficial if a lot of **duplicate data** could be eliminated before the computation in the reduce step. ✓ **more precisely?** 0,75

d) An RDD transformation is an operation that can be done on each RDD **block** separately. Example of such an operation is, **map()**, **sort()**, **ReduceByKey()**. 0,75

f) collect picks out the n number of lines from an RDD and leaves the rest of the lines. This collecting a sample. No 0

g) Streaming is a continuous stream of data. Can be used to calculate distance of moving object for example.

In Spark:



AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN 1

Blad nummer: Sheet number:
7

8 ① cluster\_point = point.map(lambda p: randint(0,1), (1, p))  
 ② total\_cluster = cluster\_point.reduceByKey(lambda a, b: a[0] + b[0],  
 a[1] + b[1])  
 average\_clust = total\_cluster.map(lambda x: x[0], (x[1][1] / x[1][0]))  
 average\_clust.broadcast()  
 ③ for i in range(1)  
 ④ clusterpoint = point.map(lambda p: closest\_clust(p), (1, p))  
 ⑤ total\_cluster = cluster\_point.reduceByKey(lambda a, b:  
 a[0] + b[0], a[1] + b[1])  
 average\_clust = total\_cluster.map(lambda x: x[0], (x[1][1] / x[1][0]))  
 ⑥ def closest\_clust(p)  
 if (distance(p, average\_clust[0]) >= distance(p, average\_clust[1]))  
 return 1  
 else  
 return 0

AID-nummer: AID-number:	1287	Datum: Date:	1/11-19
Utbildningskod: Education code:	TDDE31	Modul: Module:	TEN 1

Blad nummer: Sheet number:
8

10

assuming:

all-points is all current points

the key is the class

the value is the point

point is the point to classify

k is the number of neighbours we want to check.

Then the code would look like:

point-distance = all-points.map(lambda x: 0, (x[0], distance(point, x[1])).

sorted-distance = point-distance.sortBy(lambda x: x[1][1])

closest-points = sorted-distance.take(k)

total-closest = closest-points.reduceByKey(lambda a,b:a[0]+b[0])

new-point-class = total-closest[0][0] / k

3

new-point-class does now have the class for which the point should be according to k-nearest neighbours.