⬅ **Google Discussions**

**Exam Professional Data Engineer All Questions**
View all questions & answers for the Professional Data Engineer exam

**Go to Exam**

📄 **EXAM PROFESSIONAL DATA ENGINEER TOPIC 1 QUESTION 87 DISCUSSION**

Actual exam question from Google's Professional Data Engineer
Question #: 87
Topic #: 1
[All Professional Data Engineer Questions]

You've migrated a Hadoop job from an on-prem cluster to dataproc and GCS. Your Spark job is a complicated analytical workload that consists of many shuffling operations and initial data are parquet files (on average 200-400 MB size each). You see some degradation in performance after the migration to Dataproc, so you'd like to optimize for it. You need to keep in mind that your organization is very cost-sensitive, so you'd like to continue using Dataproc on preemptibles (with 2 non-preemptible workers only) for this workload.
What should you do?

A. Increase the size of your parquet files to ensure them to be 1 GB minimum.

B. Switch to TFRecords formats (appr. 200MB per file) instead of parquet files.

C. Switch from HDDs to SSDs, copy initial data from GCS to HDFS, run the Spark job and copy results back to GCS.

D. Switch from HDDs to SSDs, override the preemptible VMs configuration to increase the boot disk size.

**Show Suggested Answer**

by 👤 **madhu1171** at *March 14, 2020, 2:44 p.m.*

## Comments

Type your comment...

**Submit**

**rickywck** `Highly Voted` 5 years, 1 month ago

Should be A:

https://stackoverflow.com/questions/42918663/is-it-better-to-have-one-large-parquet-file-or-lots-of-smaller-parquet-files
https://www.dremio.com/tuning-parquet/

C & D will improve performance but need to pay more $$

👍 ↩ 🏳 upvoted 70 times

---

**grshankar9** 3 months, 2 weeks ago

Switching to SSDs definitely increases the cost, eliminating C & D.

👍 ↩ 🏳 upvoted 1 times

---

**diluvio** 3 years, 7 months ago

It is A . please read the links above

👍 ↩ 🏳 upvoted 5 times

---

**odacir** 2 years, 4 months ago

https://cloud.google.com/architecture/hadoop/migrating-apache-spark-jobs-to-cloud-dataproc#optimize_performance

👍 ↩ 🏳 upvoted 1 times

---

**raf2121** 3 years, 9 months ago

Point for discussion - Another reason why it can't be C or D.
SSD's are not available on pre-emptible Worker nodes (answers didn't say whether they wanted to switch from HDD to SDD for Master nodes)
https://cloud.google.com/architecture/hadoop/hadoop-gcp-migration-jobs

👍 ↩ 🏳 upvoted 8 times

---

> **rr4444** 2 years, 10 months ago
>
> You can have local SSDs for the dataproc normal or preemptible VMs
> https://cloud.google.com/dataproc/docs/concepts/compute/dataproc-pd-ssd
>
> 👍 ↩ 🏳 upvoted 1 times

---

> **raf2121** 3 years, 9 months ago
>
> Also for Shuffling Operations, one need to override the preemptible VMs configuration to increase boot disk size. (Second half of answer D is correct but first half is wrong)
>
> 👍 ↩ 🏳 upvoted 1 times

**Load full discussion…**

---

**madhu1171** `Highly Voted` 5 years, 1 month ago

Answer should be D

👍 ↩ 🏳 upvoted 12 times

---

**jvg637** 5 years, 1 month ago

D: # By default, preemptible node disk sizes are limited to 100GB or the size of the non-preemptible node disk sizes, whichever is smaller. However you can override the default preemptible disk size to any requested size. Since the majority of our cluster is using preemptible nodes, the size of the disk used for caching operations will see a noticeable performance improvement using a larger disk. Also, SSD's will perform better than HDD. This will increase costs slightly, but is the best option available while maintaining costs.

👍 ↩ 🏳 upvoted 15 times

---

> **ch3n6** 4 years, 10 months ago
>
> C is correct. D is wrong. they are using 'dataproc and GCS', not related to boot disk at all .
>
> 👍 ↩ 🏳 upvoted 3 times

---

> > **VishalB** 4 years, 9 months ago
> >
> > C is recommended only -
> > If you have many small files, consider copying files for processing to the local HDFS and then copying the results back
> >
> > 👍 ↩ 🏳 upvoted 1 times

---

> > > **FARR** 4 years, 8 months ago
> > >
> > > File sizes are already within the expected range for GCS (128MB-1GB) so not C.
> > > D seems most feasible
> > >
> > > 👍 ↩ 🏳 upvoted 3 times

---

**rajshiv** `Most Recent` 2 weeks, 2 days ago

Selected Answer: C

C is correct. It cannot be D as increasing boot disk size does not impact shuffle performance much. We need local SSDs specifically attached for shuffle storage (temporary fast storage), not just a bigger persistent boot disk.

C is the correct answer because SSD + HDFS shuffle layer = fastest for Spark shuffle-heavy jobs, while still using preemptibles and keeping costs down. The job as mentioned is shuffle-intensive. In Spark, shuffling (moving data between nodes) is heavily disk I/O intensive. Faster local storage (i.e., SSDs) can dramatically speed up shuffle operations compared to using standard HDDs. GCS is great for object storage, not shuffle storage.

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **rajshiv** 2 weeks, 2 days ago

Selected Answer: C

C is the correct answer

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **oussama7** 1 month, 2 weeks ago

Selected Answer: C

Improves shuffle management by using HDFS instead of GCS.
SSDs speed up access to temporary data.
Compatible with Dataproc's preemptible cost model, without requiring more non-preemptible workers.

👍 ↩ 🏳 upvoted 2 times

⊟ 👤 **Parandhaman_Margan** 1 month, 3 weeks ago

Selected Answer: C

Preemptible Cost Considerations

Using preemptibles (with 2 non-preemptible workers) is cost-effective, but shuffle operations still need fast local storage. SSDs improve reliability without increasing instance costs significantly

👍 ↩ 🏳 upvoted 2 times

⊟ 👤 **f74ca0c** 3 months, 4 weeks ago

Selected Answer: A

A.
Not D because it doesn't make sens to move to SSD when cost-senstive

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **Javakidson** 6 months ago

A is the answer

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **SamuelTsch** 6 months, 2 weeks ago

Selected Answer: A

I think either A or C. The problem is occured by I/O performance. Option A is feasible, which reduces the number of files leading better parallel processing. Option C tries to handle I/O performance issue.
Taking other factors like budget and no mention of HDD/SSD, option A is possible the correct answer.

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **baimus** 7 months, 2 weeks ago

Selected Answer: A

There's no mention of a drive type used, only GCS. That means A is the only sensible option.

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **987af6b** 9 months, 2 weeks ago

Selected Answer: A

Question doesn't actually say they are using HDD in the scenario, for that reason I choose A

👍 ↩ 🏳 upvoted 2 times

⊟ 👤 **philli1011** 1 year, 2 months ago

A
We don't know if HDD was used, so we can know what to do about that, but we know that the parquet files are small and much, and we can act on that by increasing the sizes to have lesser number of it.

👍 ↩ 🏳 upvoted 2 times

⊟ 👤 **rocky48** 1 year, 5 months ago

Selected Answer: A

Should be A:
https://stackoverflow.com/questions/42918663/is-it-better-to-have-one-large-parquet-file-or-lots-of-smaller-parquet-files

👍 ↩ 🏳 upvoted 1 times

⊟ 👤 **rocky48** 1 year, 5 months ago

Given the scenario and the cost-sensitive nature of your organization, the best option would be:

C. Switch from HDDs to SSDs, copy initial data from GCS to HDFS, run the Spark job, and copy results back to GCS.

Option C allows you to leverage the benefits of SSDs and HDFS while minimizing costs by continuing to use Dataproc on preemptible VMs. This approach optimizes both performance and cost-effectiveness for your analytical workload on Google Cloud.

👍 ↩ 🚩 upvoted 1 times

---

👤 **Mathew106** 1 year, 9 months ago

Selected Answer: A

https://stackoverflow.com/questions/42918663/is-it-better-to-have-one-large-parquet-file-or-lots-of-smaller-parquet-files

Cost effective is the key in the question.

👍 ↩ 🚩 upvoted 1 times

---

👤 **Nandhu95** 2 years, 1 month ago

Selected Answer: D

Preemptible VMs can't be used for HDFS storage.
As a default, preemptible VMs are created with a smaller boot disk size, and you might want to override this configuration if you are running shuffle-heavy workloads.

👍 ↩ 🚩 upvoted 1 times

---

👤 **midgoo** 2 years, 2 months ago

Selected Answer: D

Should NOT be A as:
1. The file size is already at the optimal size
2. If the current file size works well in the current Hadoop, it is expected to have similar performance in Dataproc

The only difference between the current and Dataproc is that Dataproc is using preemptible nodes. So yes, it may incur a bit more cost by using SSD but assuming using the preemptible already save most of it, so we want to save less to improve the performance

👍 ↩ 🚩 upvoted 1 times

> 👤 **Mathew106** 1 year, 9 months ago
>
> Optimal size is 1GB
>
> 👍 ↩ 🚩 upvoted 1 times

---

👤 **[Removed]** 2 years, 2 months ago

Selected Answer: A

Cost sensitive is the keyword.

👍 ↩ 🚩 upvoted 1 times

**Load full discussion...**