➢ **Title:**

Cloud Vs. On-Prem for Generative AI System

➢ **Source:**

https://youtu.be/r7cdRV0fDho?si=d-XSme3ffsh73c079

➢ **Objective:**

To understand the difference between using cloud platform and on-premises infrastructure for running generative AI systems and there impact on cost, security and performance.

➢ **What Concept Covered in Video**

✓ **Cost:**

- Cloud is cheaper in the beginning, but long-term use can be expensive

- On-Prem is need more money in the beginning but can be cheaper later.

✓ **Security:**

- Cloud Depends on outside companies, also still they have good security.

- On-Prem gives full control over your data.

✓ **Speed and Performance:**

- Cloud depends on internet and providers server speed

- On-prem might give better speed for some tasks.

✓ **Updates:**

- Cloud gets the latest updates and new tools quickly.

- On-Prem needs manual updates and changes.

✓ **Scalability:**

- Cloud can be increase or decrease storage.

- On-Prem System cannot be changed easily and have fixed capacity.

➢ **What I Learned**

From this video, l learned that both cloud and on-prem systems have good and bad sides

✓ Cloud is good for flexibility and new tools.

✓ On-prem is better when you need more control over data and cost.

✓ Sometimes companies use both like mix of cloud and on-Prem this is called as hybrid approach.

➢ **Screenshots:**