

LiCO H100 Cluster Kickstarter

เริ่มต้นใช้งาน H100 Slurm Cluster บน LiCO ชั้นเบื้องต้น

มีเวอร์ชันภาษาอังกฤษ



ณภัทร ศรีจันทร์

15 มีนาคม 2569

วิทยาลัยการคอมพิวเตอร์ มหาวิทยาลัยขอนแก่น

หัวข้อ

1. เกริ่นนำ
2. ความรู้พื้นฐาน
3. เขียน Hello World บน LiCO
4. เทรนโมเดลบน LiCO
5. ใช้งานการ์ดจอบน LiCO
6. เทรนโมเดลโดยไม่ล่ม
7. สรุป

เกริ่นนำ



ทำไม

ทำไม

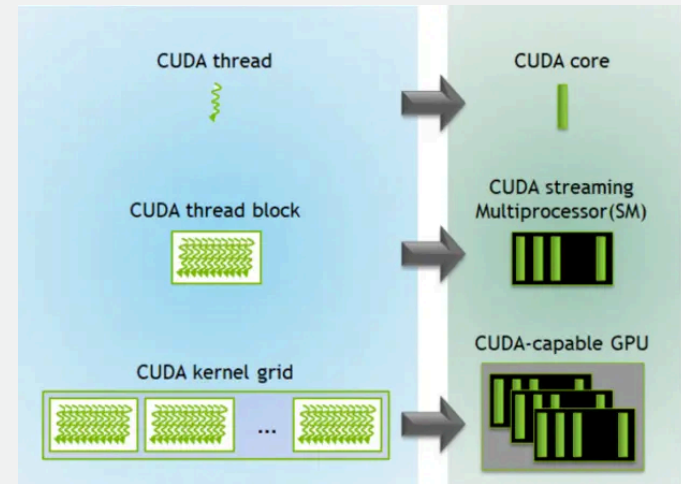
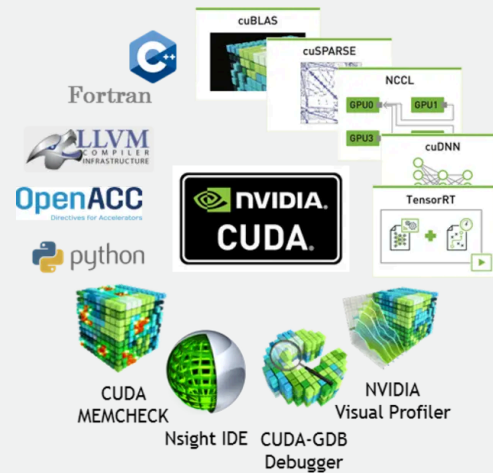
เครื่องมือสำคัญ

ทำไม

เครื่องมือสำคัญ , ฟังพาตนเองได้

CUDA

Compute Unified Device Architecture¹² .

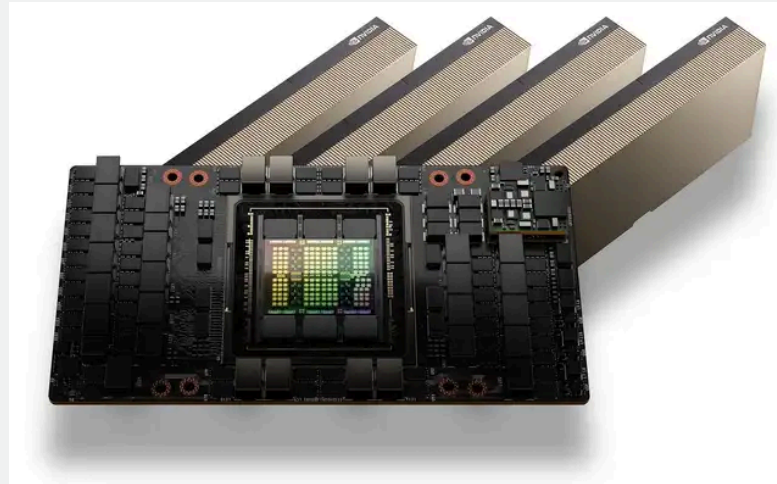


¹<https://developer.nvidia.com/cuda>

²<https://www.modular.com/blog/democratizing-ai-compute-part-3-how-did-cuda-succeed>

H100

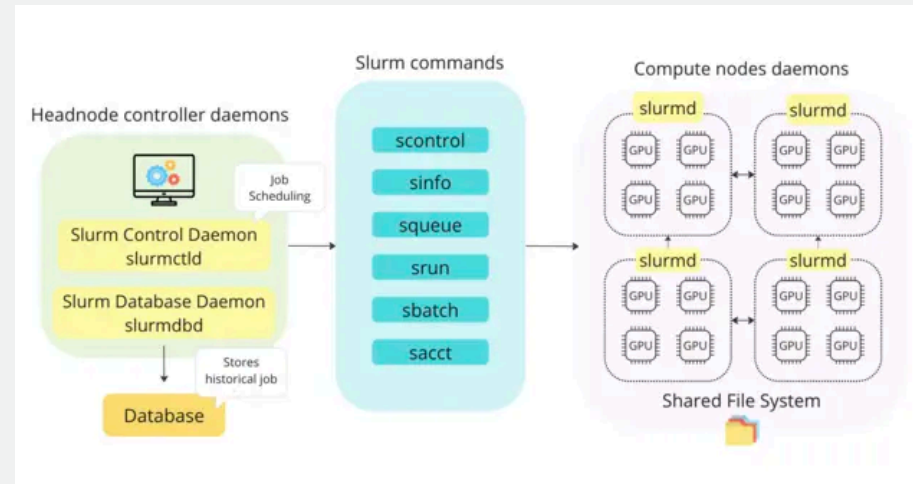
การ์ดจอระดับอุตสาหกรรมจาก NVIDIA (รุ่น Hopper)³



³<https://www.nvidia.com/en-us/data-center/h100/>

Slurm

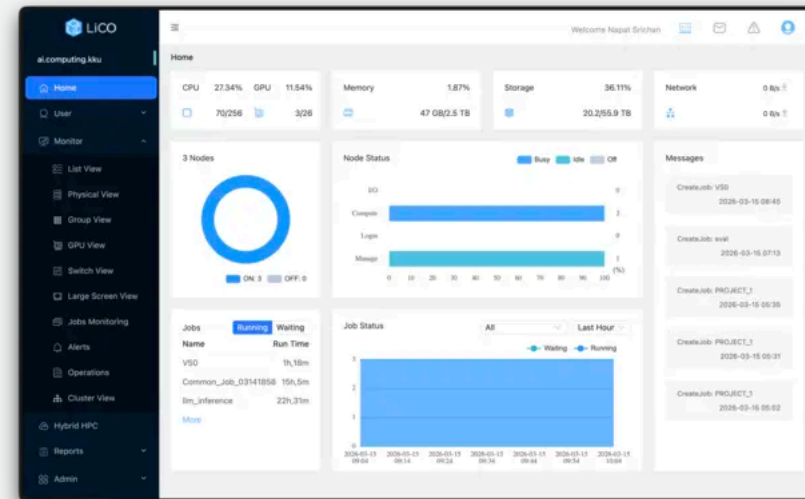
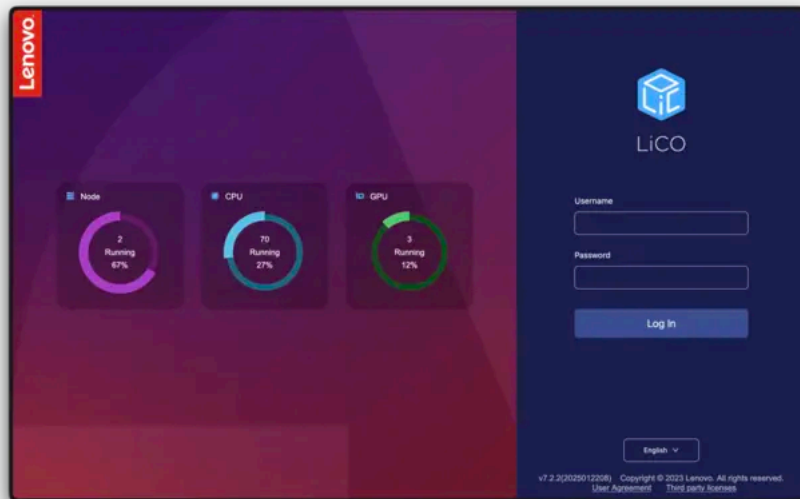
โปรแกรมจัดการภาระงาน/ทรัพยากรสำหรับเซิร์ฟเวอร์หลายเครื่อง⁴



⁴<https://slurm.schedmd.com/overview.html#architecture>

LiCO

โปรแกรม Slurm โดย Lenovo ควบคุมผ่านเว็บได้



Kaen-Coon (แก่นคูณ)



Summarized flow

1. ล็อกอินเข้าระบบ LiCO
2. ใส่โปรแกรมและข้อมูลลงไป
2. สร้าง Job
3. รัน Job และตรวจสอบสถานะ
4. ดึงผลลัพธ์ออกมา

เงื่อนไขที่ควรมี

- โปรแกรมเขียนโค้ด: VSCode หรือโปรแกรมอื่น ๆ เทียบเท่า
- ชำนาญภาษา Python และเคยใช้ PyTorch มาก่อน
- รู้จักคำสั่ง command line จะดีมาก; สามารถแก้ bug จาก logs ได้บ้าง
- เข้าถึงเครือข่ายภายใน มข. ได้ (ต้อง VPN เข้ามาหากอยู่ภายนอก)
- Username/password (ยกมือขอ staff หากยังไม่ได้)

เขียน/ส่งโปรแกรมเข้า LiCO กัน!

LiCO ที่ <https://10.198.253.15:8000>

https://github.com/anonymawew/lico-training/tree/main/01-hello_world

ข้อสังเกต

- Slurm คือการเขียนโปรแกรม shell
-
-

ข้อสังเกต

- Slurm คือการเขียนโปรแกรม shell
- เราเห็นการคิวงานทำงานจริง
-

ข้อสังเกต

- Slurm คือการเขียนโปรแกรม shell
- เราเห็นการคิวงานทำงานจริง
- เวอร์ชัน Python ค่อนข้างเก่า (3.6)

ลองเทรนโมเดลบน LiCO

https://github.com/anonymaw/lico-training/tree/main/02-mnist_trainnig

ข้อสังเกต

- ต้องใช้คอนเทนต์เนอร์เฉพาะ
-
-

ข้อสังเกต

- ต้องใช้คอนเทนเนอร์เฉพาะ
- จำเป็นต้องใช้ virtual environment
-

ข้อสังเกต

- ต้องใช้คอนเทนเนอร์เฉพาะ
- จำเป็นต้องใช้ virtual environment
- ค่อนข้างช้า (เรายังใช้ CPU อยู่)

ใช้การ์ดจอเทรนโมเดลบน LiCO

https://github.com/anonymawew/lico-training/tree/main/03-mnist_trainnig_cuda

ข้อสังเกต

- เทรนเร็วขึ้นมากด้วย GPU
-
-

ข้อสังเกต

- เทรนเร็วขึ้นมากด้วย GPU
- ต้องระบุว่าใช้ MIG ทุกครั้ง
-

ข้อสังเกต

- เทรนเร็วขึ้นมากด้วย GPU
- ต้องระบุว่าใช้ MIG ทุกครั้ง
- ความคืบหน้าหายถ้ากดยกเลิก/เกิดข้อผิดพลาด

คะแนนพิเศษ: เทรนโมเดลโดยใช้ checkpoint

(ใช้ ChatGPT ได้) เขียน 2 methods:

โหลดโมเดลจาก checkpoint ถ้ามี, บันทึกรหัสโมเดลทุก ๆ training epoch

สรุปจบ

- LiCO → Slurm → Script → GPU
- PyTorch เท่านั้น CUDA คือกุญแจสำคัญ
- โปรดใช้ในขณะที่ยังมีให้ใช้ฟรี
- User จะถูกลบในวันพรุ่งนี้
- สงสัยถามได้