1. Goal: Demonstrate and characterize GPU power variations in ML parallel workloads. This could enable power-aware time/space-sharing and placement performance optimizations
   a. **Profiling**: Measurable metrics that affect variability at device level
      i. Jetson:
         a) SGEMM metrics with/without DVFS. Measurements with fine-grained frequency control is a good approximation to understand range of variations
         b) Fine-grained sensor measurements contrasted with nvprof reporting
      ii. V100:
         a) Testbench: Automated toolchain install, runs and reporting
         b) Concurrent SGEMM runs across 4xV100 (space locality), Multiple jobs on a single GPU (to represent time locality). Extend this to future measurements
      iii. Obtain power, frequency, temperature metrics across all GPUs using nvprof
   b. **Characterization**: Create stressmark suite to demonstrate variability by picking workloads that span applications/bottlenecks representing realistic usage in space-sharing systems
      i. SGEMM/DGEMM
      ii. RESNET (Language)
      iii. BERT (Vision)
      iv. DLRM (Recommender systems)
   c. **Mitigation**: Decided to defer to future work post project-proposal.
2. Challenges
   a. C4130 nodes require force reboot after 30m or so with CUDA installation. This seems to be a known issue, but the provided workaround isn't working yet. This limits long running simulations [Mailing list link](#)
   b. SGEMM kernel fails for matrix dimensions larger than 16k x 16k
   c. More than 4 GPU variation study likely not possible. Even if we get two nodes of c4130, the relative placement in cluster is unknown to have reproducible effects
3. Timeline

| Week starting | Goals |
|---|---|
| April 10 | • Co-locate cross combinations of the following from the [Nvidia DL repository](#) with SGEMM (base vs test)<br>a) RESNET (Vision)<br>b) BERT (Language)<br>c) DLRM (Recommender) |
| April 17 | • Collect metrics for same scenarios with GPU-boost disabled<br>• Collate database with metrics captured across all runs include 1-4 GPU combinations<br>• Formalize variation using model properties and trends/anomalies from profiling |
| April 24 | Analysis, documentation, and presentation |

4. Resource request from course staff
   a. We have c4130 node reserved until the 20$^{th}$. It might be helpful if we can get an extension on this for a week or 2 nodes of c4130.
   b. Guidance on analysis plan since the data collection part feels open ended