



A Lightweight GPU Monitoring Extension for Pegasus Kickstart

George Papadimitriou, Ewa Deelman

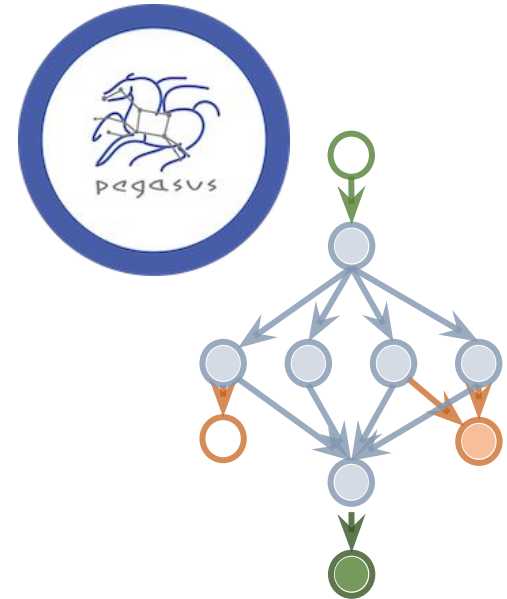
Workshop on Workflows in Support of Large-Scale Science (WORKS)
November 15th, 2021

This work was funded by the US Department of Energy under contract **DE-SC0012636M** and by NSF under contract **1664162**.

Pegasus Kickstart: What is it?

Pegasus Kickstart (or Kickstast for short)

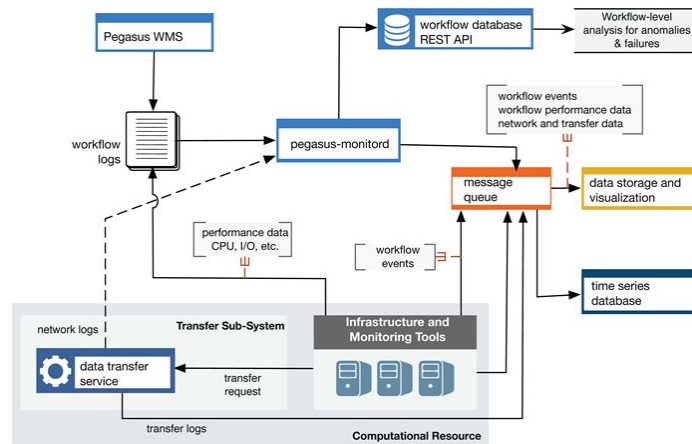
- A lightweight job wrapper
- Implemented in C
- Captures performance and provenance data
- YAML output
- Part of the Pegasus workflow management system tools



Kickstart in Panorama: What's new?

Pegasus Panorama extends Kickstart

- Online collection of performance statistics (traces) as captured by Linux's procs
- Correlates them with workflow jobs
- Interval as fast as 1 second
- AMQP publish capability
- JSON output



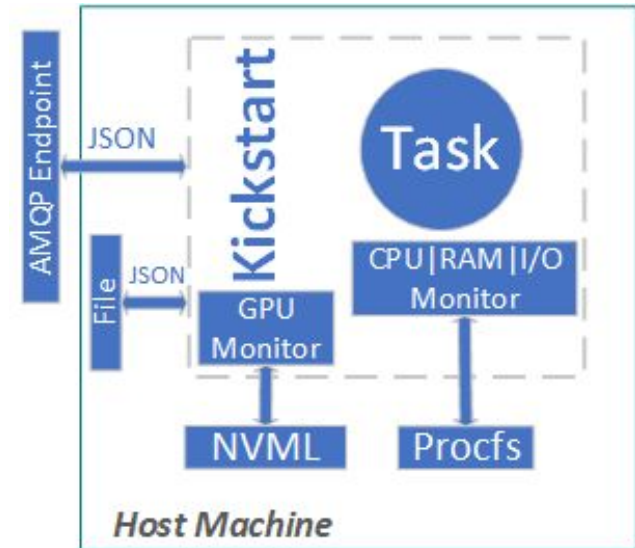
<https://panorama360.github.io>

Now with Nvidia Monitoring Support !

Kickstart: Nvidia Monitoring - Approach

- Leverages Nvidia's Monitoring Library (NVML)
 - C-based API
 - Interfacing with the Nvidia driver
- Kickstart implements a lightweight C wrapper for NVML
- Kickstart polls for new GPU statistics as fast as 1 second
- Multi-threading support

<https://docs.nvidia.com/deploy/nvml-api/index.html>



Kickstart online monitoring

Kickstart: Nvidia Monitoring - Events

- **kickstart.inv.online.gpu.env:** Contains information about the GPU environment - produced once at the beginning of the job
- **kickstart.inv.online.gpu.stats:** Contains a snapshot of the GPU counters - produced throughout the execution
- **kickstart.inv.online.gpu.stats.max:** Summary of max values observed - produced after the job finishes

kickstart.inv.online.gpu.env

Field	Type	Description
ts	int	seconds since epoch of when the event was created
hostname	string	hostname of the compute node
vf_label	string	pegasus workflow label
nvidia_driver_version	string	MAJOR.MINOR driver of the execution environment
site	string	pegasus execution site
information	string	pegasus transformation
task_id	string	pegasus task id
cuda_version	float	MAJOR.MINOR cuda version available on the execution environment
dag_job_id	string	pegasus dag job id
vf_job_id	string	pegasus workflow uuid
condor_job_id	string	condor job id
devices_count	int	number of detected gpu devices
devices_compute_node	string	device compute node as by admin
devices_id	int	gpu device id in the environment
devices_is_cuda_capable	bool	has cuda capability
devices_max_clock	int (MHz)	max streaming multiprocessor clock speed
devices_max_mem_clock	int (MHz)	max gpu memory clock speed
devices_power_limit	int (kilowatt)	gpu power limit
devices_total_memory	int (bytes)	total gpu memory
devices_total_bar1_memory	int (bytes)	total shared memory allocated for 3rd party devices
devices_name	string	gpu product name
devices_cuda_capability	float	MAJOR.MINOR
devices_max_gpu_clock	int (MHz)	max gpu clock
devices_max_temperature	int (Celsius)	max temperature
devices_max_gpu_utilization	int (%)	max graphics core utilization

kickstart.inv.online.gpu.stats.max

Field	Type	Description
ts	int	seconds since epoch of when the event was created
hostname	string	hostname of the compute node
vf_label	string	pegasus workflow label
site	string	pegasus execution site
information	string	pegasus transformation
task_id	string	pegasus task id
dag_job_id	string	pegasus dag job id
vf_job_id	string	pegasus workflow uuid
condor_job_id	string	condor job id
devices_id	int	gpu device id in the environment
devices_max_bar1_mem_usage	int (bytes)	maximum memory used for communication with 3rd party devices
devices_max_temp	int (Celsius)	maximum observed temperature
devices_max_gpu_utilization	int (%)	max graphics core utilization
devices_product_name	string	gpu product name
devices_power_usage	int (kilowatt)	max power usage
devices_gpu_memory_usage	int (bytes)	max gpu memory usage
devices_unique_bus_id	string	unique bus id
ts	int	seconds since epoch of when the event was created
hostname	string	hostname of the compute node
vf_label	string	pegasus workflow label
site	string	pegasus execution site
information	string	pegasus transformation
task_id	string	pegasus task id
dag_job_id	string	pegasus dag job id
vf_job_id	string	pegasus workflow uuid
condor_job_id	string	condor job id
id	int	gpu device id in the environment
compute_node	string	device compute node as by admin
name	string	gpu product name
bus_id	string	unique bus id
pci	string	pci address
vendor_name	int (hex)	vendor name in hex
device_name	int (hex)	device name in hex
gpu_clock	int (MHz)	graphics clock speed
gpu_mem_clock	int (MHz)	shared multiprocessor clock speed
mem_clock	int (MHz)	memory clock speed
power_usage	int (kilowatt)	power usage
mem_usage	int (bytes)	gpu memory usage
mem_utilization	int (%)	gpu memory utilization - percent of mem speed to memory use since last clock change
gpu_utilization	int (%)	gpu utilization - percent of mem speed to gpu use
bar1_mem_usage	int (bytes)	memory usage for communication with 3rd party devices
temperature	float	maximum temp observed in the environment
pci_is_exposed	int (boolean)	boolean indicated on the pci bus during 400ns sample period
pci_is_collected	int (boolean)	boolean on the pci bus during 400ns sample period
compute_gpu_mem_usage	int (bytes)	memory usage of the compute node
compute_gpu_mem_usage_max	int (bytes)	memory usage of the compute node
graphics_gpu_mem_usage	int (bytes)	graphics memory usage of the graphics processor
graphics_gpu_mem_usage_max	int (bytes)	graphics memory usage of the graphics processor
gpu_utilization_max	int (%)	gpu utilization
gpu_mem_utilization_max	int (%)	gpu mem utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization
mem_utilization_max	int (%)	memory utilization

Kickstart: Nvidia Monitoring - Levels

- **DEFAULT:** Reports general status of the GPU and statistics of the CUDA kernels
- **MONITORING_PCIE:** Enables monitoring of the pcie bus for data transmission to and from the GPU
- **GRAPHICS_PROCS:** Captures statistics of the graphics (video) processes

Monitoring Level	Measurement Time
DEFAULT	~6ms
MONITORING_PCIE	~228ms

More information on GitHub

<https://github.com/pegasus-isi/pegasus/tree/panorama/src/tools/pegasus-kickstart/nvidia>



Kickstart: Nvidia Monitoring - Invocation

- **KICKSTART_MON_URL:** Sets the location the statistics will be saved (AMQP or file)
- **-m INT:** Enables online monitoring with interval INT
- **-G:** Enables default gpu monitoring
- If you want to use with Pegasus consider referencing
 - <https://pegasus.isi.edu/documentation/manpages/pegasus-kickstart.html?highlight=kickstart>.
 - <https://pegasus.isi.edu/presentations/2019/pegasus-office-hours-monitoring.pdf>

```
export KICKSTART_MON_URL = \  
    rabbitmq://[USERNAME:PASSWORD]@hostname[:port]/api/exchanges/[VIRTUAL_HOST]/[EXCHANGE_NAME]/publish  
or  
export KICKSTART_MON_URL = file://filename  
pegasus-kickstart <args> -G -m 10 ./exec
```

Listing 1: Example invocation of GPU monitoring



FAQ: Do I Need to Use The Entire Pegasus WMS?

- **No!**
- **Kickstart is standalone and can be invoked independently!**
 - But if you are using Pegasus you automatically benefit from it ;)
- **Prebuilt versions of the GPU enabled Kickstart available online**
- **Offered through the lightweight pegasus worker package**

<http://download.pegasus.isi.edu/pegasus/5.1.0panorama>

QUESTIONS ?



George Papadimitriou
University of Southern California
georgpap@isi.edu

<https://pegasus.isi.edu>

<https://panorama360.github.io>

<https://github.com/pegasus-isi/pegasus/tree/panorama>



U.S. DEPARTMENT OF
ENERGY

This work was funded by the
US Department of Energy
under contract **DE-SC0012636M**,
and by NSF under contract **1664162**.

