# Data Engineering with Apache Spark

Markus Dale, medale@asymmetrik.com

May 2019

- Slides: https://github.com/medale/
- Scala Spark Code Examples: https://github.com/medale/

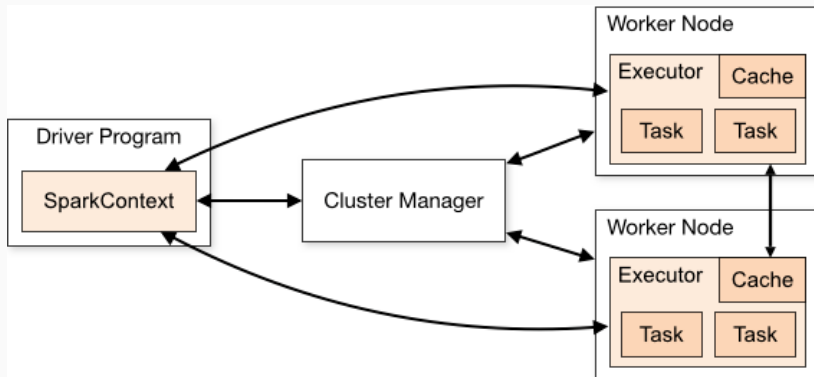Figure 1: Laptop

Figure 2: Beefed-up Server

- TODO: distributed resource management
- Spark Standalone
- Kubernetes
- Hadoop YARN
- Mesos
- Databricks
- Spark application: Driver, executors
- tasks, partitions (distributed file system, commonly accessible data store)

Source: Apache Spark website

# And now for something completely different: Colon Cancer



- Screening saves lives!
    - Colonoscopy - talk to your doc
    - Dave Barry: A journey into my colon — and yours
- Colorectal Cancer Alliance

- medale@asymmetrik.com
- Infrequent blog/past presentations http://uebercomputing.com/