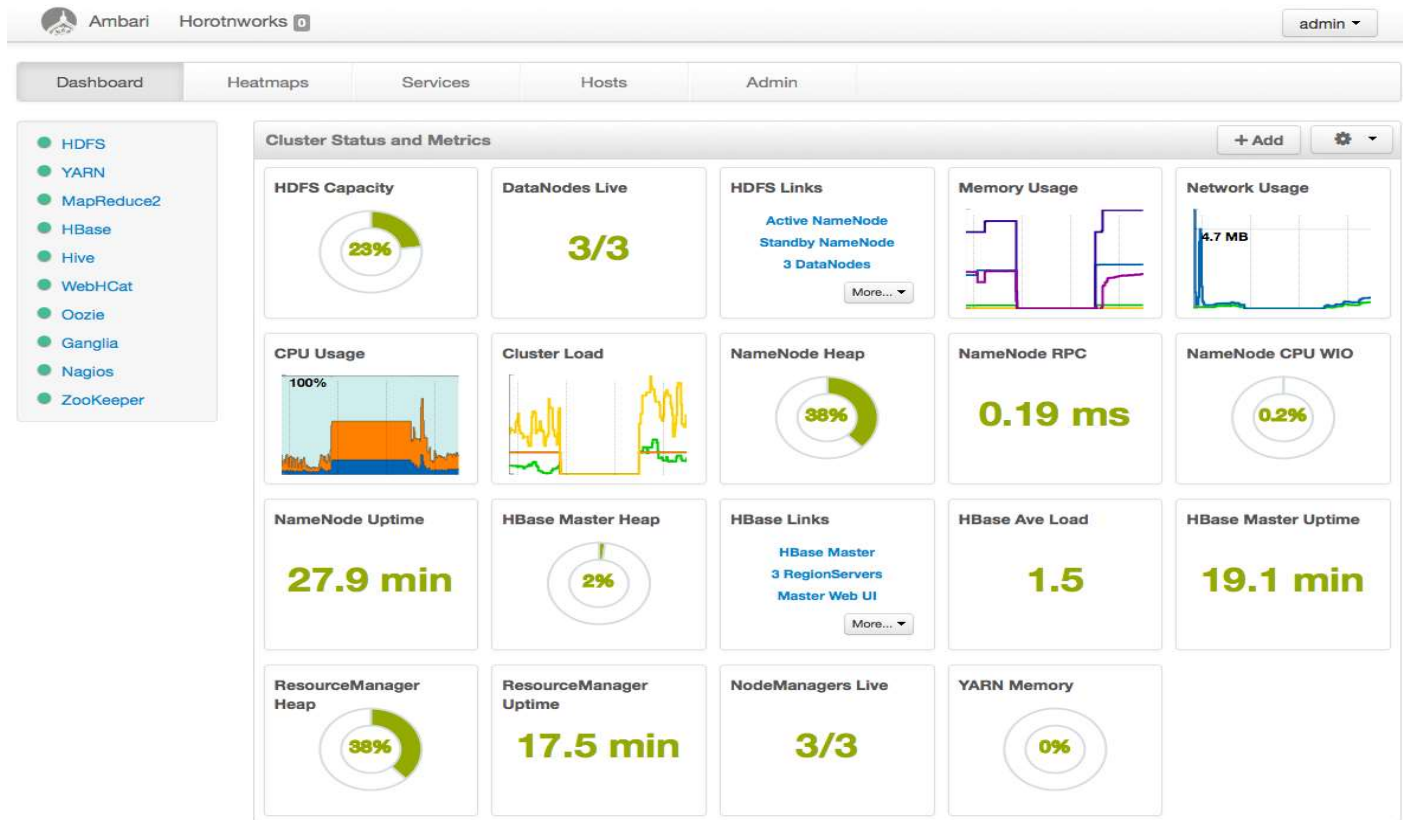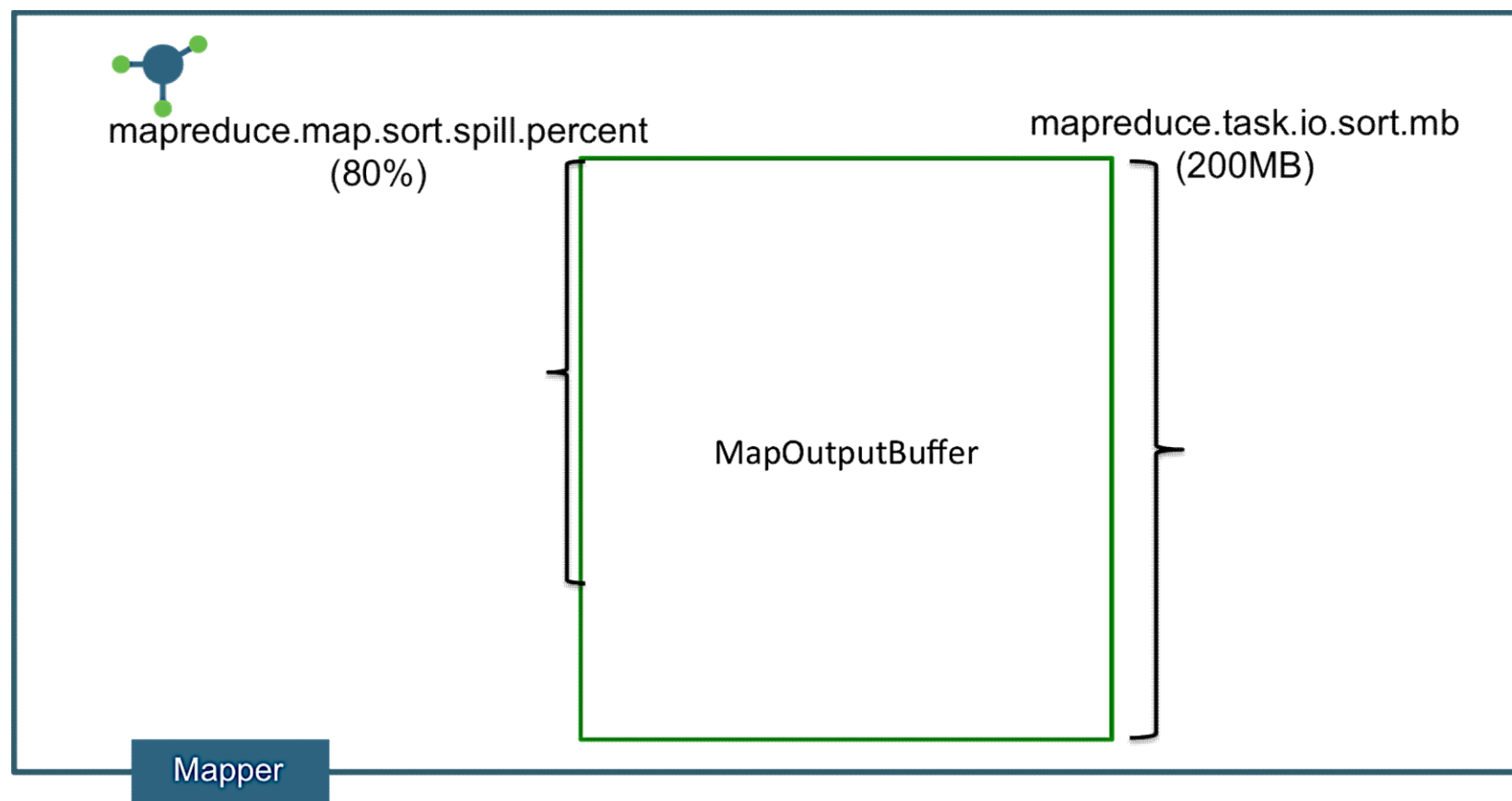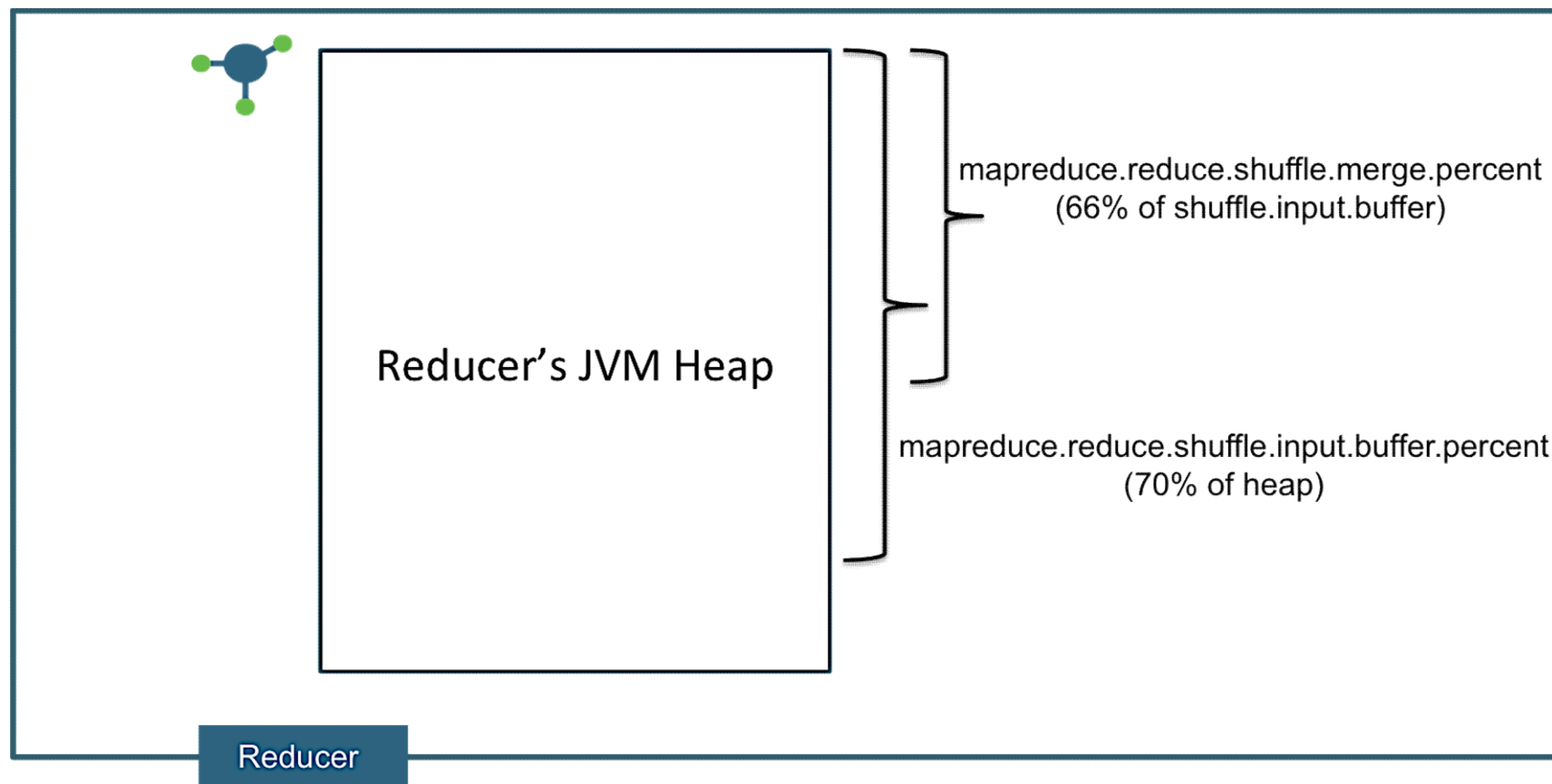# Optimizing MapReduce Jobs

# Monitoring Performance

# Optimization Best Practices

- Distribute workload evenly across NodeManagers
- Use a combiner
- Avoid instantiating new objects, including string literals
- Prefer StringBuilder over string concatenation
- Use data compression
- Avoid converting numeric types to text
- Define and configure a RawComparator
- Prefer StringUtils.split over String.split

# Optimizing the Map Phase

mapreduce.map.sort.spill.percent
(80%)

mapreduce.task.io.sort.mb
(200MB)

MapOutputBuffer

Mapper

# Optimizing the Reduce Phase

# Data Compression

- **Snappy**: org.apache.hadoop.io.compress.SnappyCodec

- **gzip**: org.apache.hadoop.io.compress.GzipCodec

- **bzip2**: org.apache.hadoop.io.compress.BZip2Codec

- **LZO**: com.hadoop.compression.lzo.LzopCodec

- **DEFLATE**: org.apache.hadoop.io.compress.DefaultCodec

# Limitations of Compression

- **Space vs. time**: is the gain in filesystem space or smaller network traffic worth the additional time it takes to compress and decompress the data?

- **Splittable vs. Non-splittable**: a major concern in MapReduce, since large files are split up into chunks across the cluster.

# Configuring Data Compression

```
Configuration conf = job.getConfiguration();
conf.setBoolean(MRJobConfig.MAP_OUTPUT_COMPRESS,
true);
conf.setClass(MRJobConfig.MAP_OUTPUT_COMPRESS_CODEC,
                        SnappyCodec.class,
                        CompressionCodec.class);


conf.setBoolean(FileOutputFormat.COMPRESS, true);
conf.setClass(FileOutputFormat.COMPRESS_CODEC,
                        SnappyCodec.class,
                        CompressionCodec.class);
```