SEGUE

parallel R
in the cloud
two lines of code

no kidding!

SEGUE

# SEGUE

syntax...

require(segue)

myCluster <- createCluster()

contratulations. we've built a hadoop cluster!

SEGUE

more syntax...

parallel apply() on lists:

base R:
lapply( X, FUN, … )

segue:
emrlapply( myCluster, X, FUN, … )

# howzit work?

## emrlapply()

**SEGUE**

list is serialized to CSV and uploaded to S3 – streaming input file

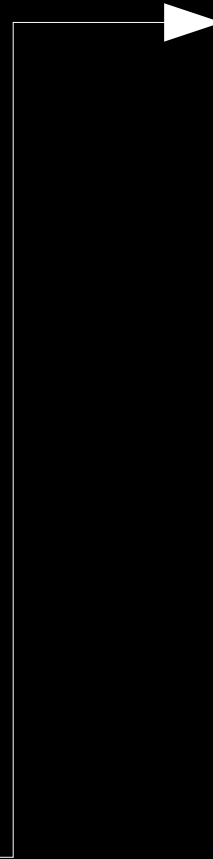function, arguments, r objects, etc are saved & uploaded

EMR copies files to nodes – mapper.R picks them up

CSV is input to mapper.R applies function to each list element

output is serialized into emr part-xxxxx files on s3

part files are downloaded to R and deserialized

deserialized results are reordered and put into a list object

# SEGUE

## when to use segue...

embarrassingly parallel

cpu bound

apply on lists with many items

object size: to / from s3 roundtrip

# SEGUE

segue project page
http://code.google.com/p/segue/

google groups
http://groups.google.com/group/segue-r