

Notes	much much slower but less memory. no use of disk	1h5 and 5 minutes for 5 test runs	started having to write to disk at 256. Fail at 2048 due to not enough disk space (30 GB allocated, probably like 13 GB free)	20 minutes for 5 test runs	Run-time isn't as dependent on nodes until it gets to larger datasets. Picture diagrams to show how work is being divided among nodes.
-------	--	-----------------------------------	---	----------------------------	--

Notes	In both cases I see a turning point if the cluster getting faster somewhere between 1 and 16, and performance improvements seem to peak at 5x for the DF and 4.3x for the RDD, even though the cluster has 6 worker nodes.
	Further questions: do more experiments with quantities between 1 and 16 to determine the "turning point" and 16 and 64 to get the increase smoothed out.
	Further questions: does the size of the cluster change where these break even points are?