# /working with larger datasets
## two-approaches
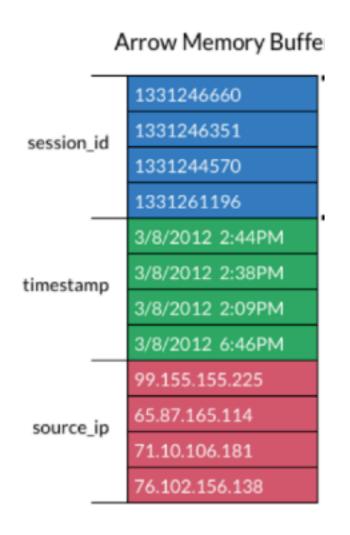
Interactivity on Large Datasets with Shiny, Arrow & Base R

# /arrow + parquet



- `{arrow}` evaluates lazily by default

- Verbs: `filter, select, mutate, join, distinct, group_by` + summarize, and across

- execution only runs on `dplyr::collect()`

- massive performance gains using parquet files and smart partitioning