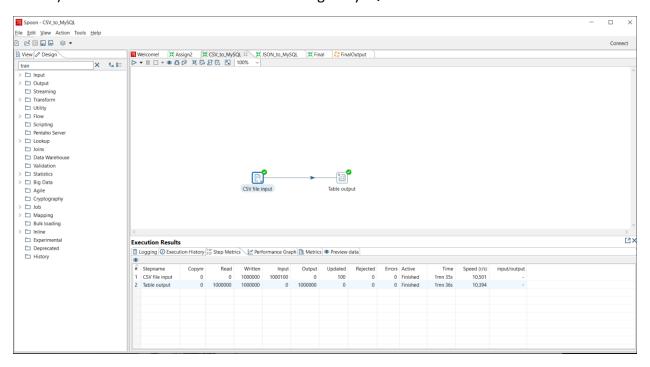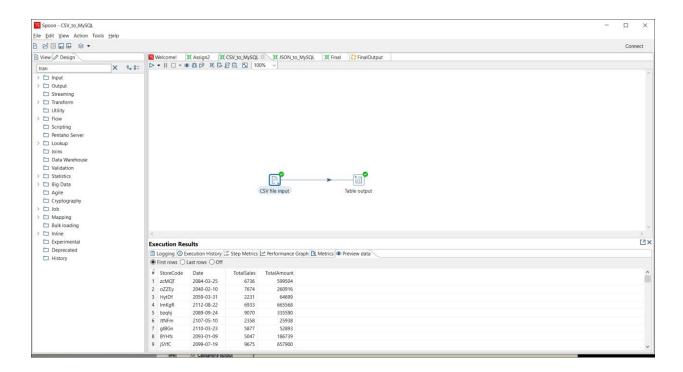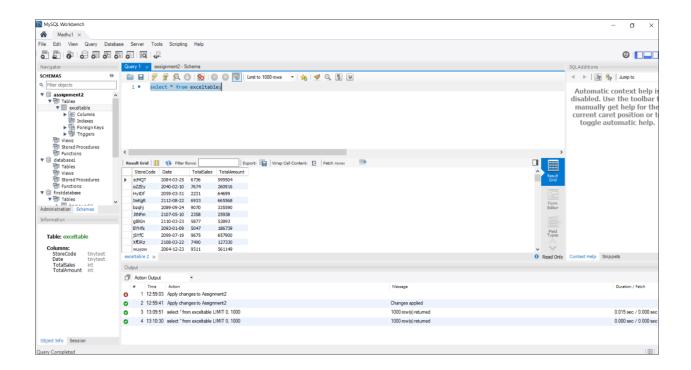Perform Extraction:
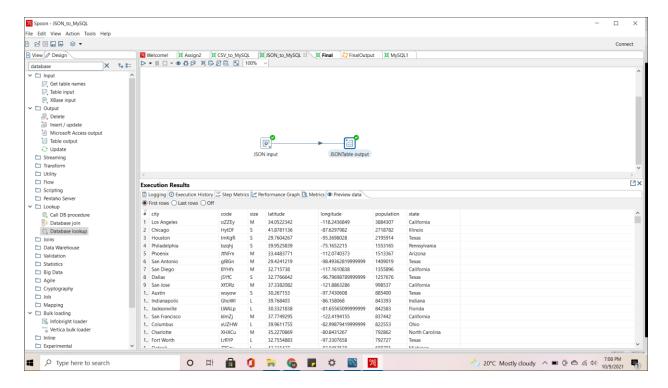
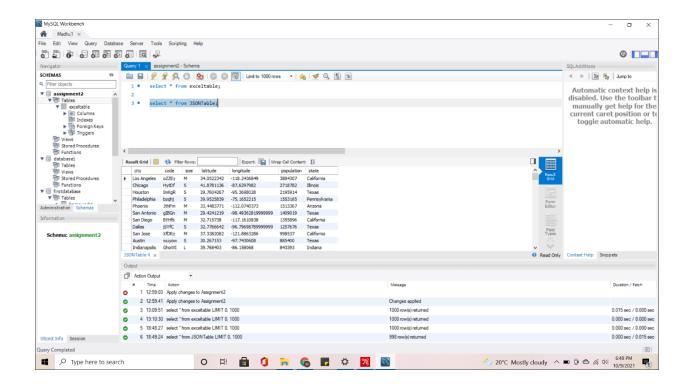2.a)      Extract all the sales CSV files to a single MySQL table.
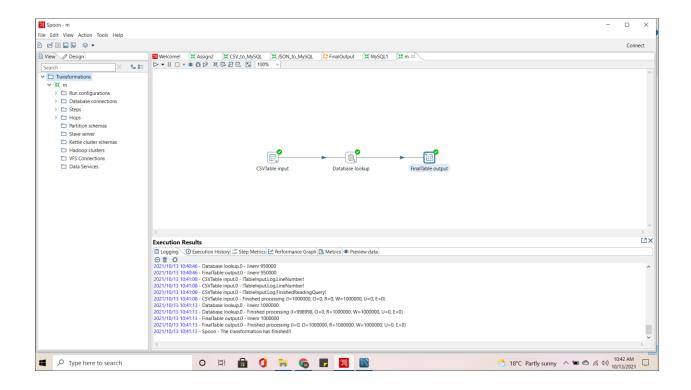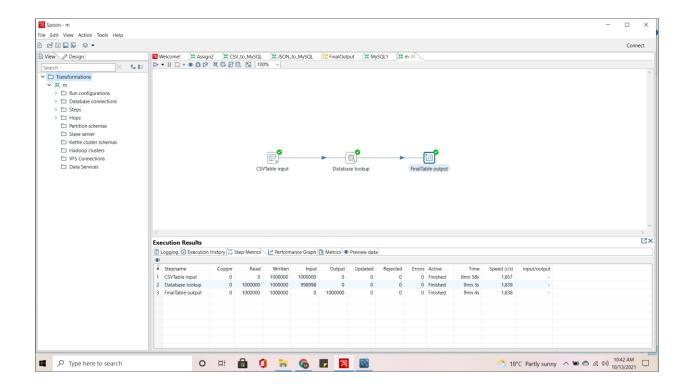
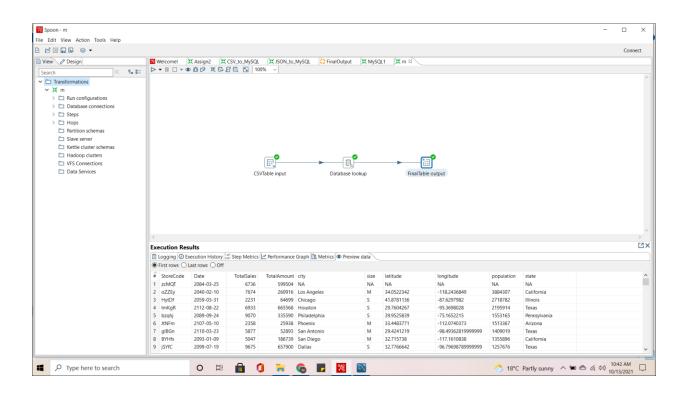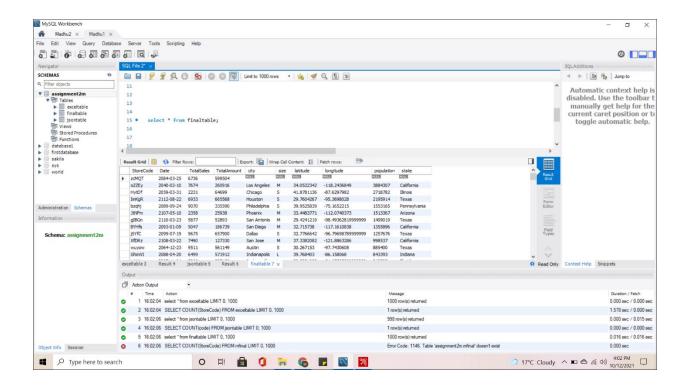## 2.b)          Extract the store JSON file to another MySQL table
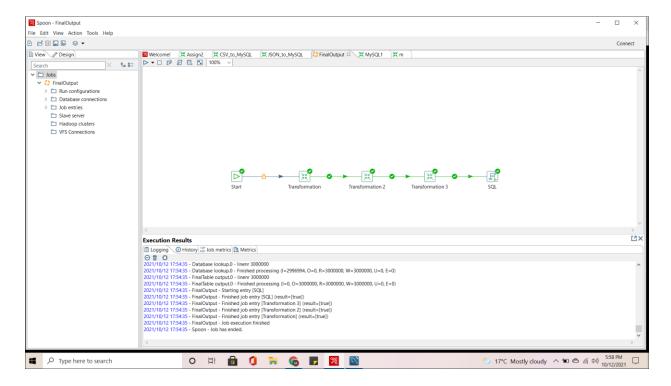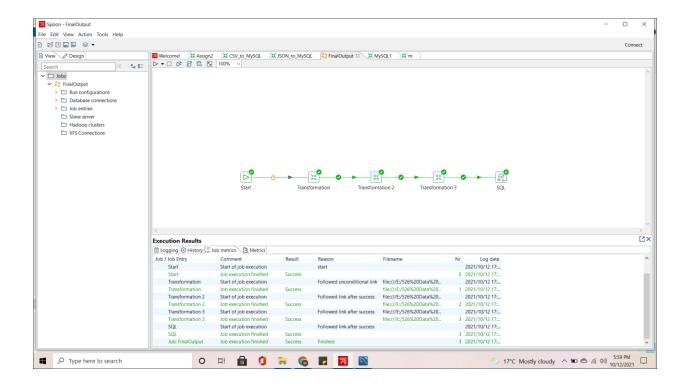
## 3.      Perform Outer Join (or Lookup):

4. Create a Kettle job to include all the transformations created above.

4.a)    For those source data rows that do not match a lookup row, set all the column values from the Lookup step to have '**NA**' as their values