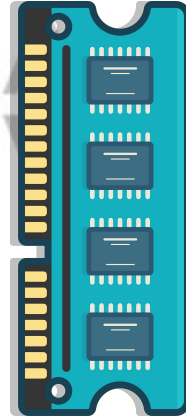
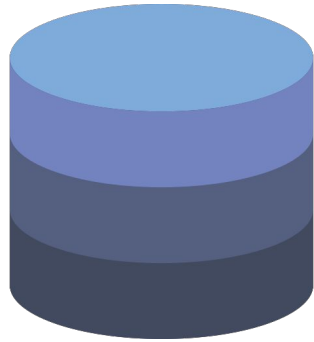


Caching and Persistence

Cache

Persist



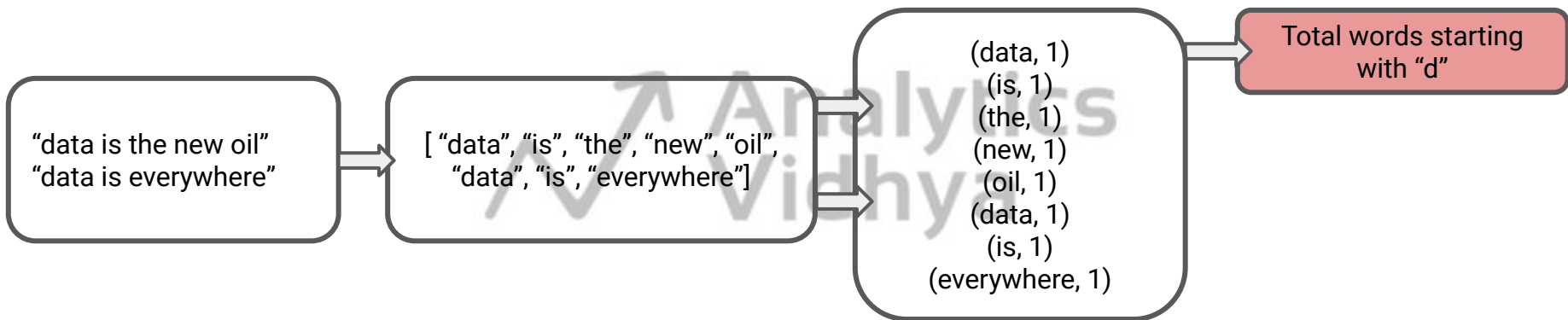
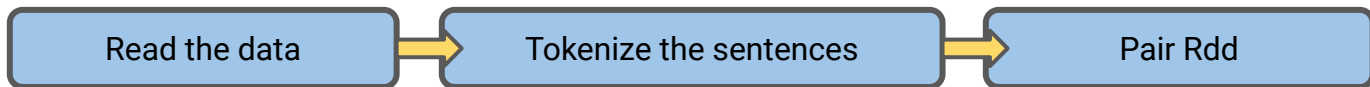
Computation

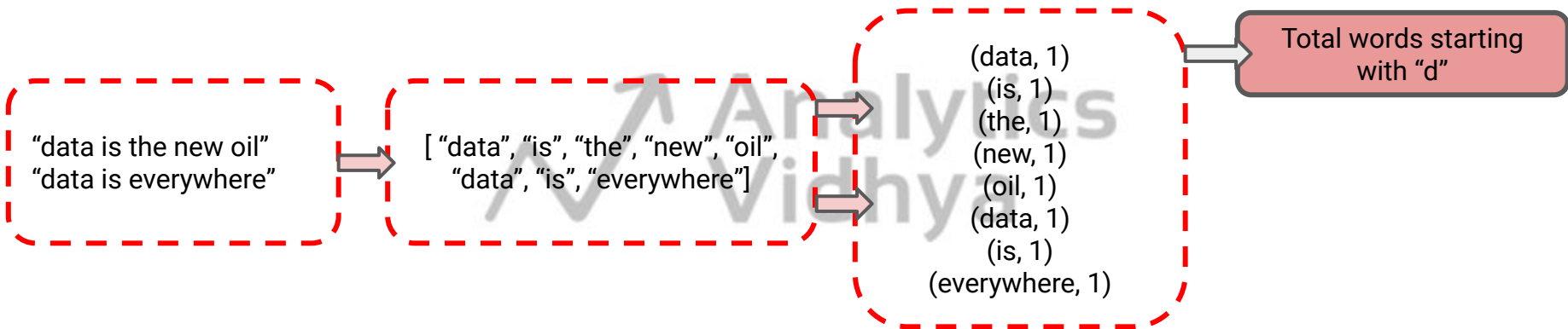
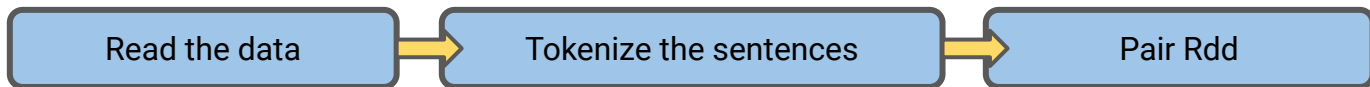


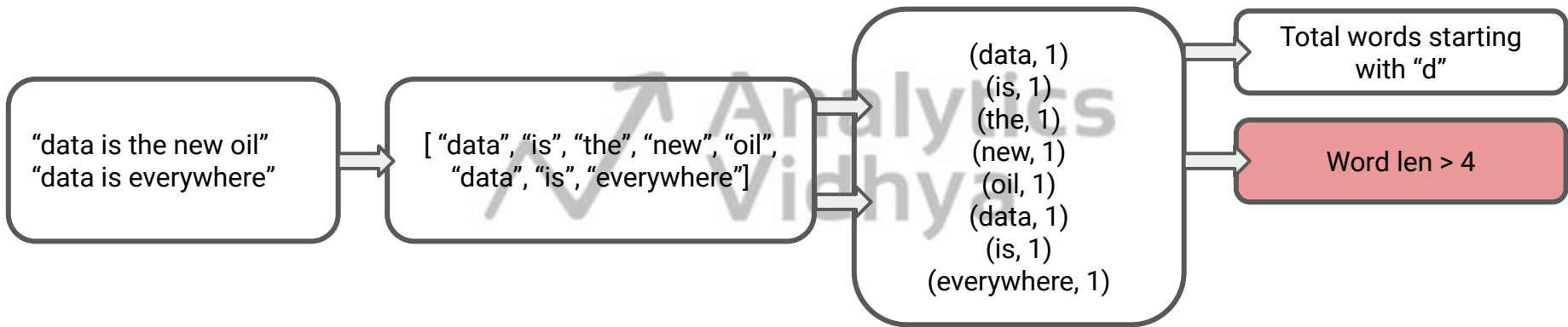
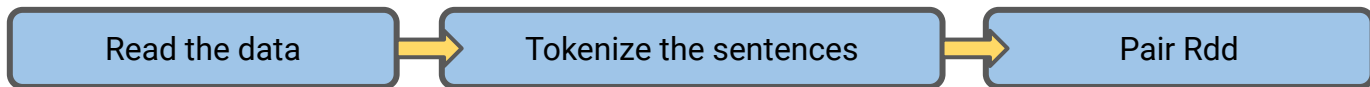
Computation

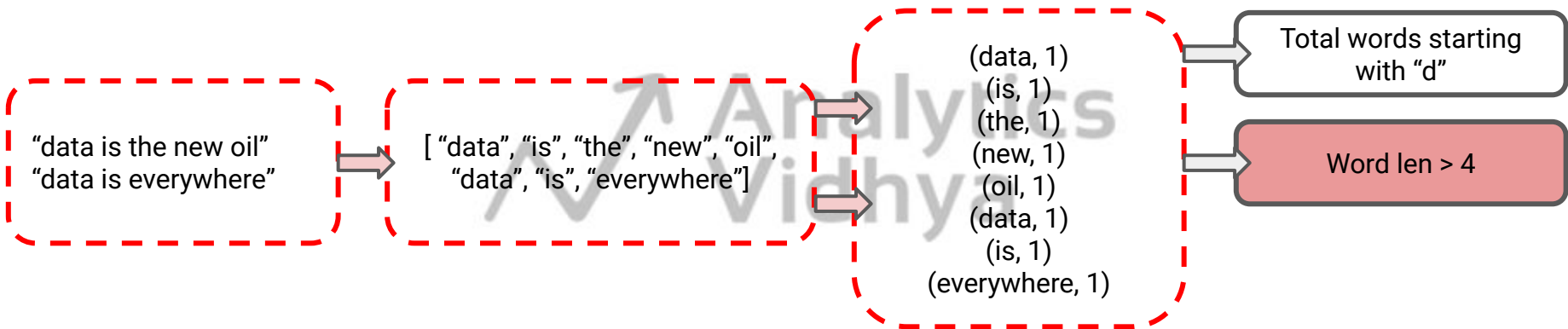
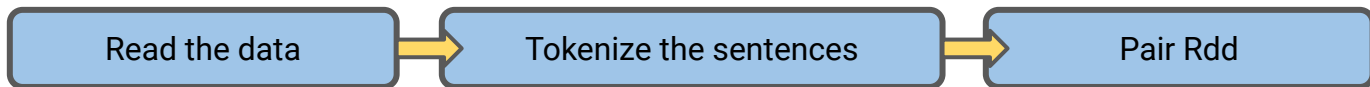


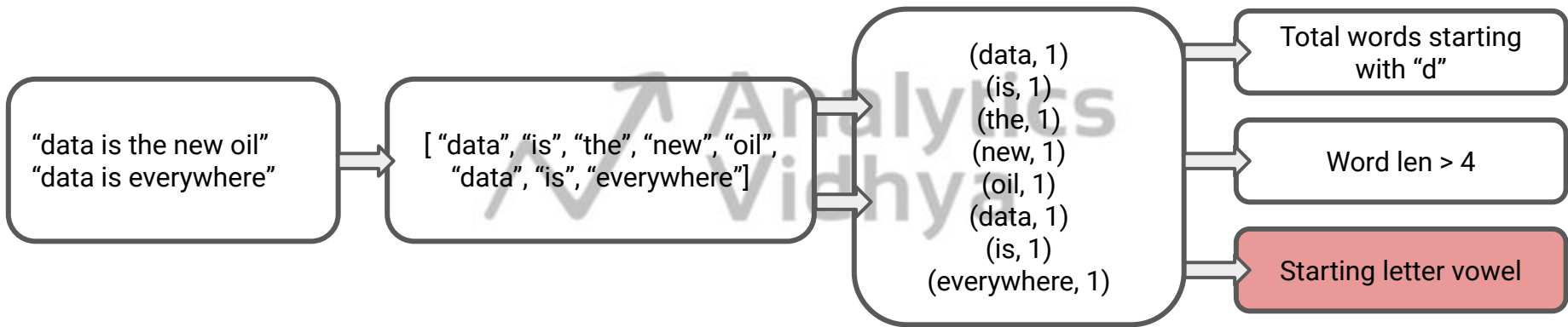
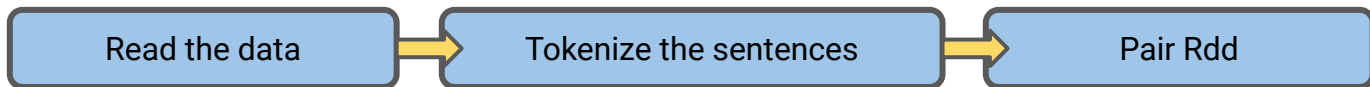
Computation

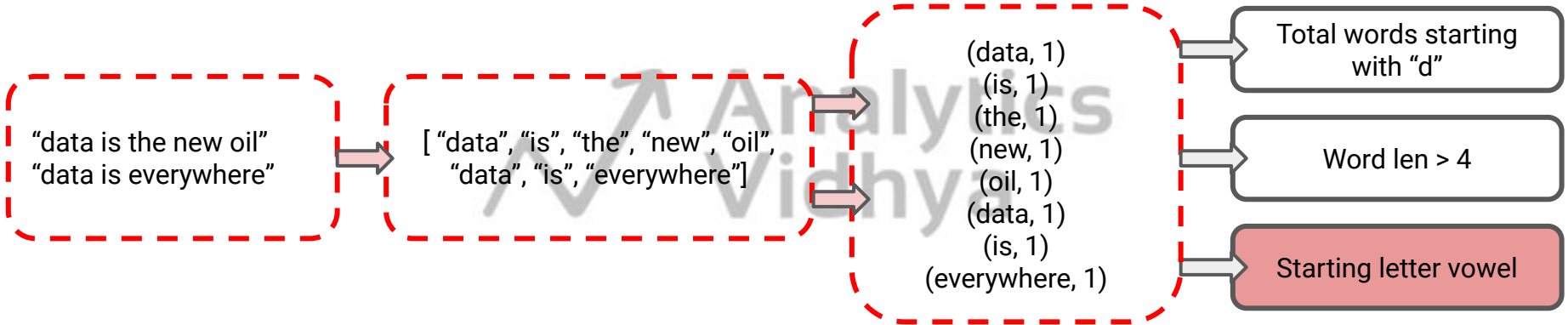
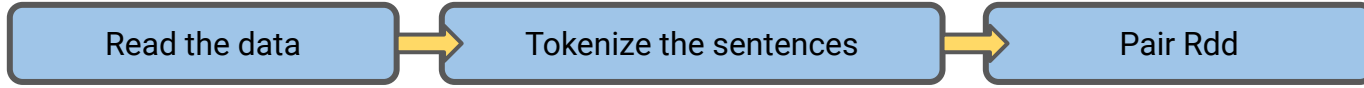


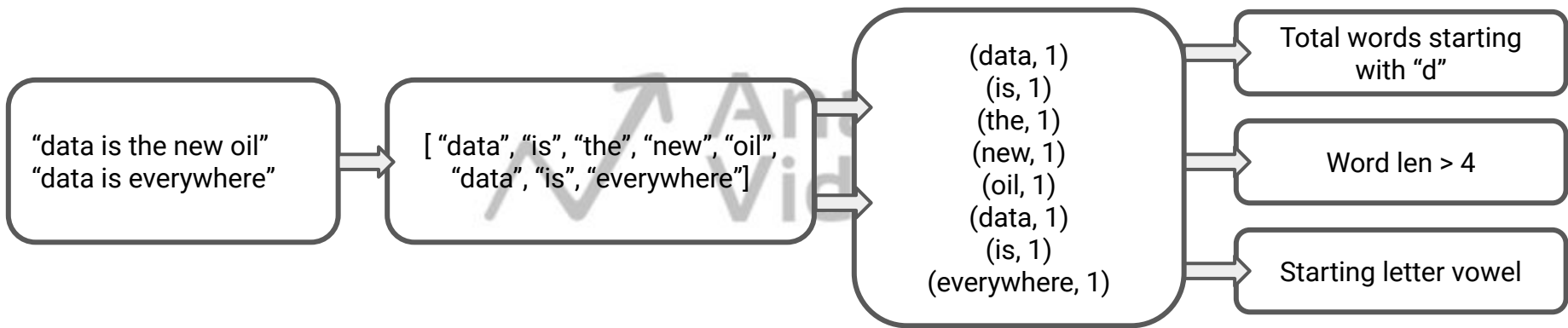
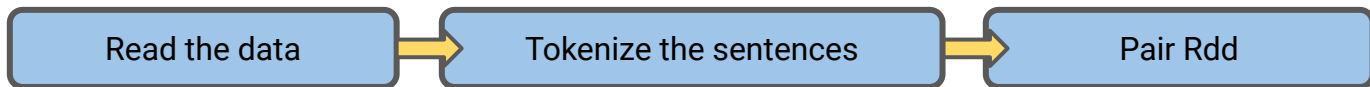




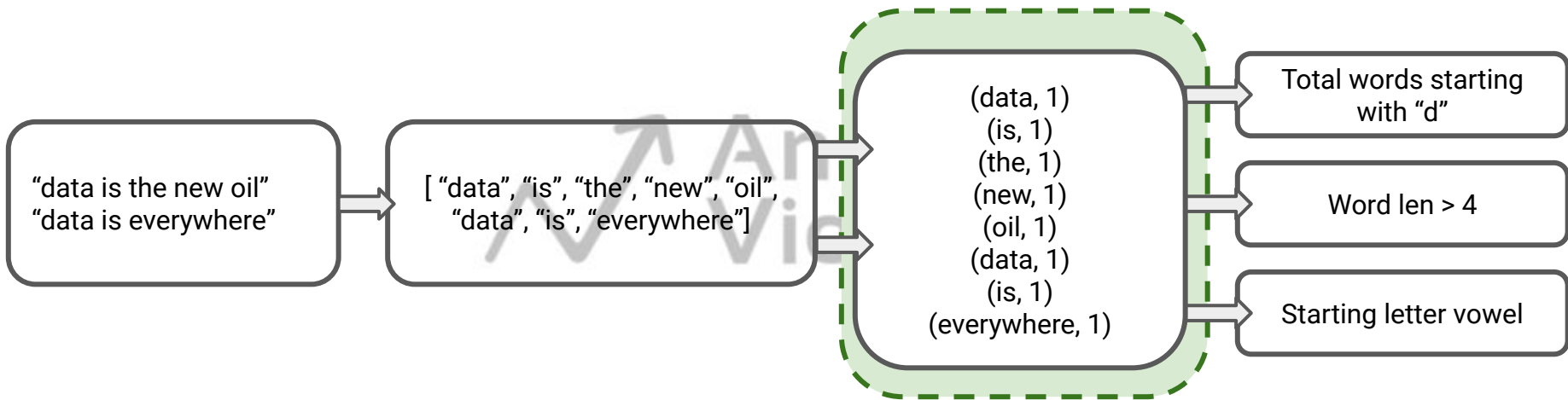
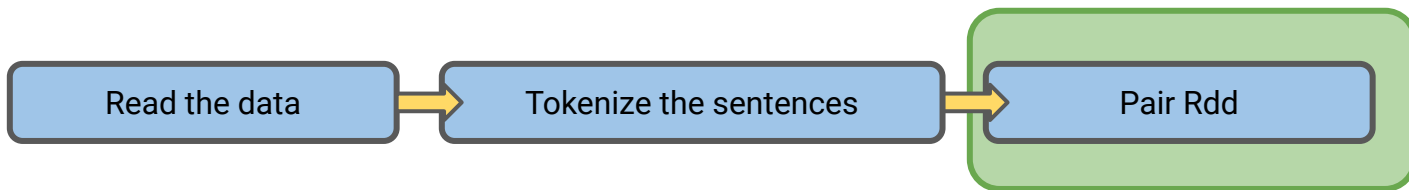




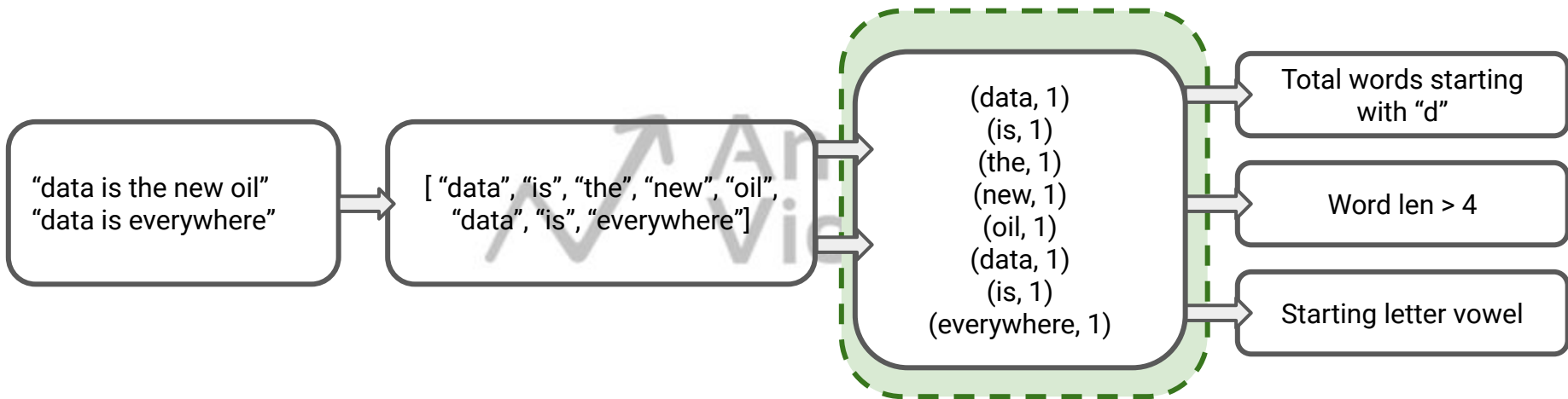
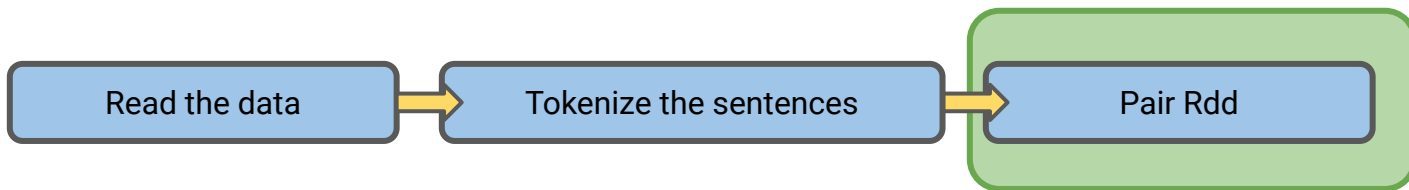


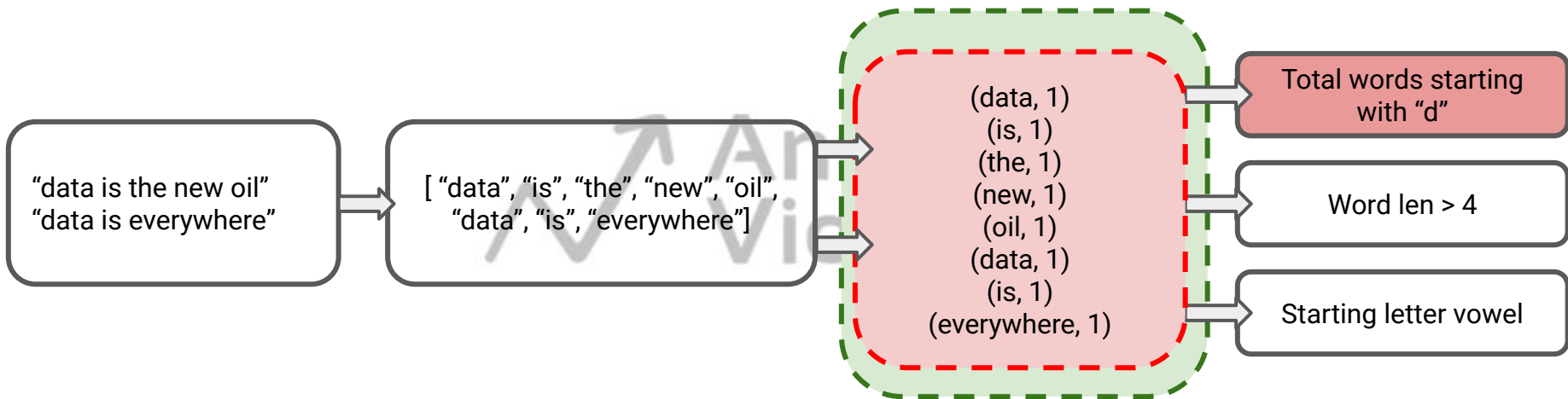
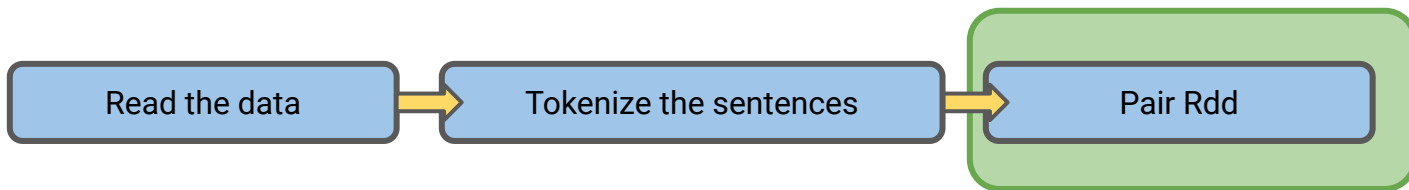


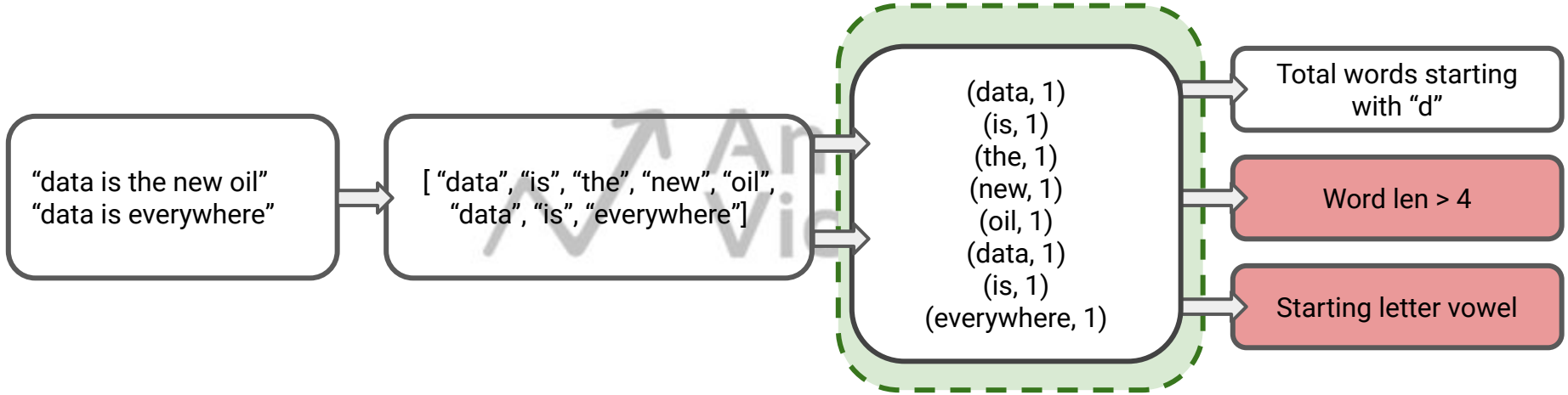
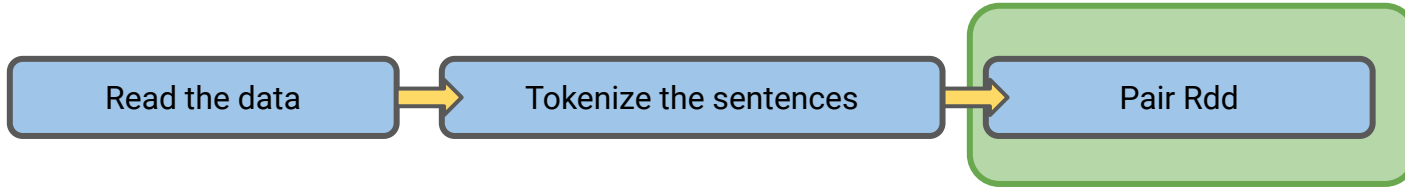
cache/persistence

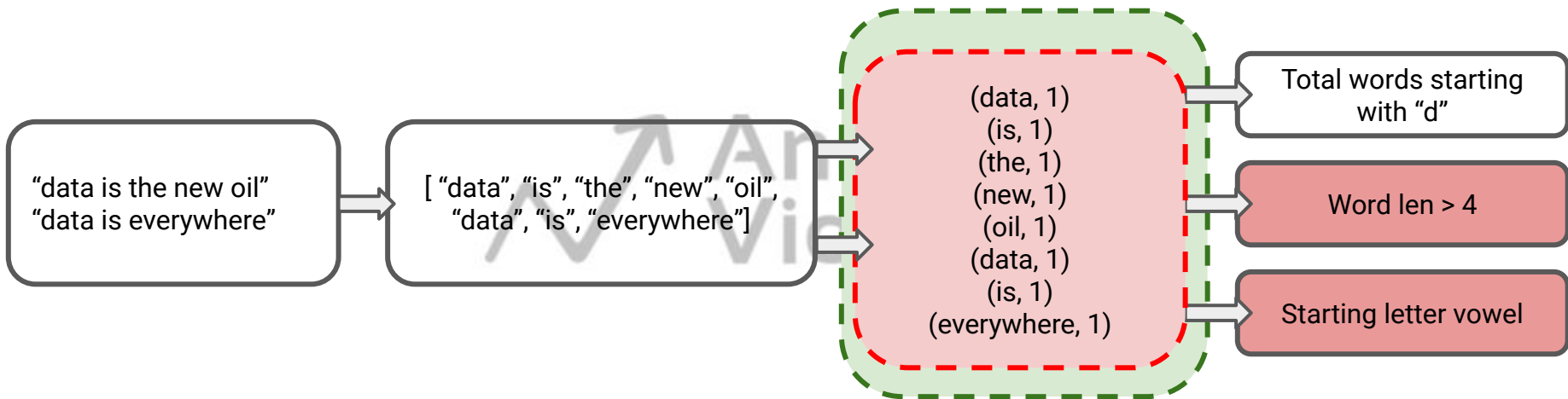
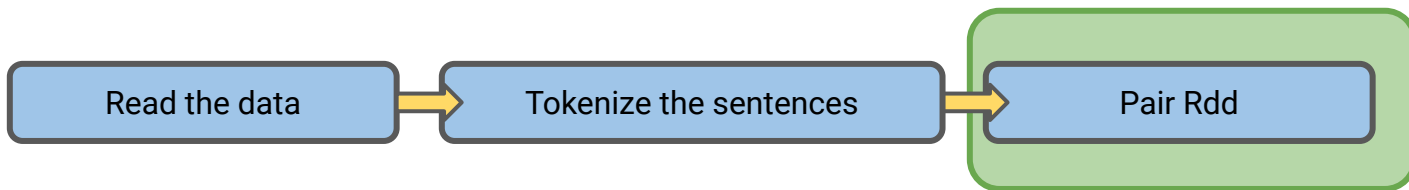


cache/persistence





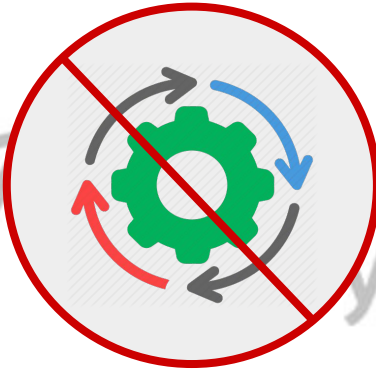




Benefits of Caching/Persistence



Saves the data by default in memory

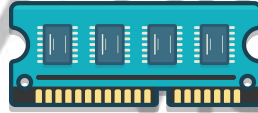
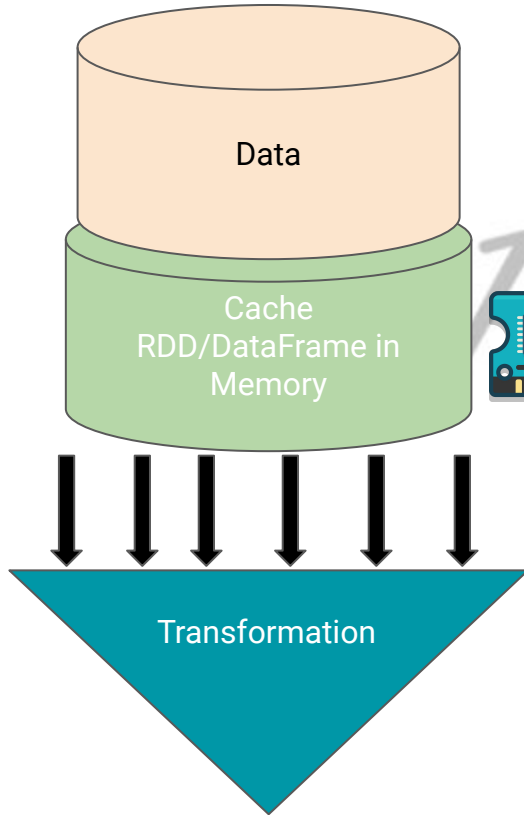


Avoids re-computation of the whole lineage



Improves Performance

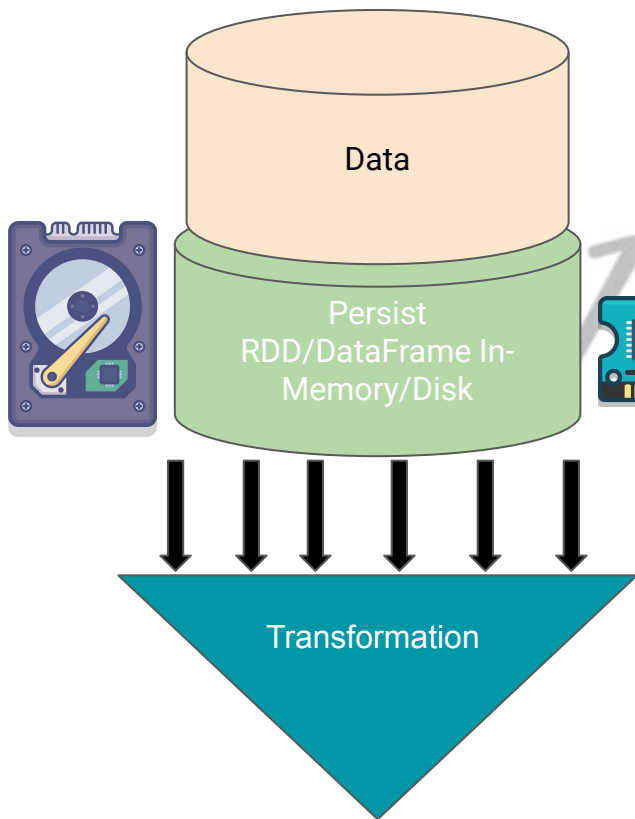
Cache



Syntax:

```
#cache RDD  
rdd.cache()
```


Persistence



Syntax:

```
#persist RDD with default storage-level  
rdd.persist()
```

```
#persist RDD with MEMORY_AND_DISK_2  
rdd.persist(pyspark.StorageLevel.MEMORY_AND_DISK_2)
```



Thank You!!