



Distributed Data Mining Lab



Meet the team



DataEngineers



Chaitanya Aggarwal



Ishmeet Kaur



Muneer Ahmad

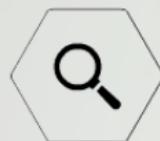


Roba Abbas

Goals



LEARN Distributed Technologies



ANALYZE Frameworks e.g Spark,Hadoop etc.



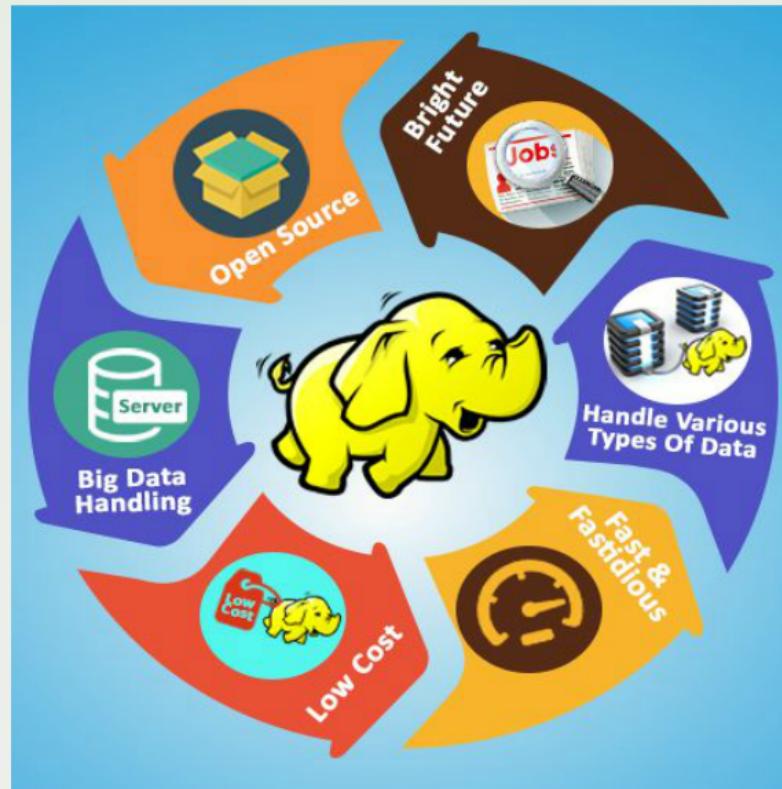
RUN real application on distributed platform

Common Scenario in companies

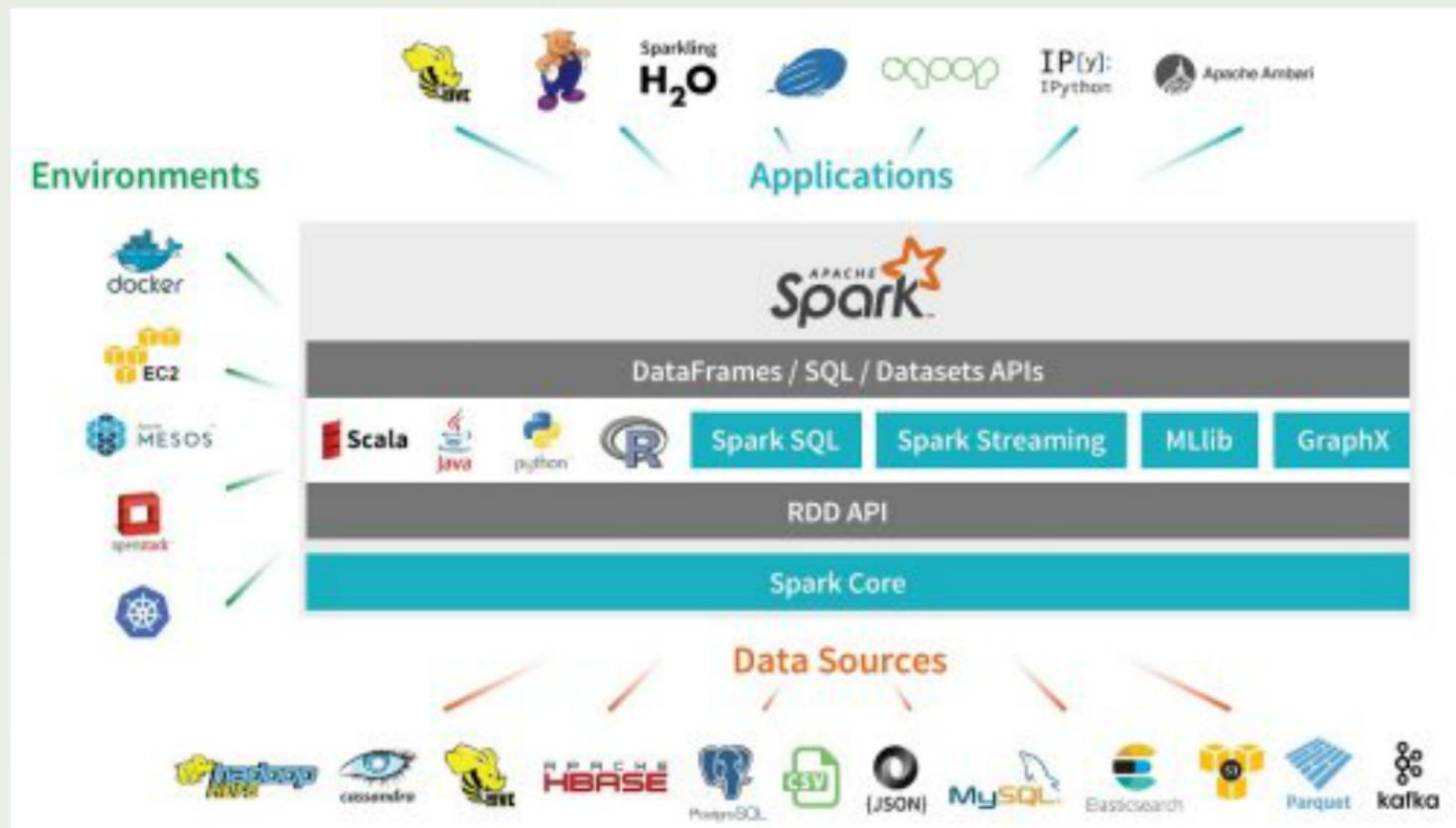




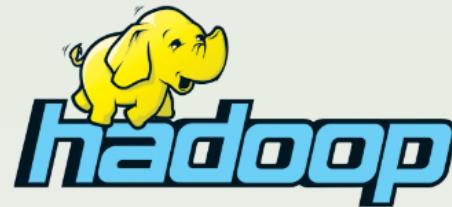
Why Hadoop?



Why Spark?



Technologies Used



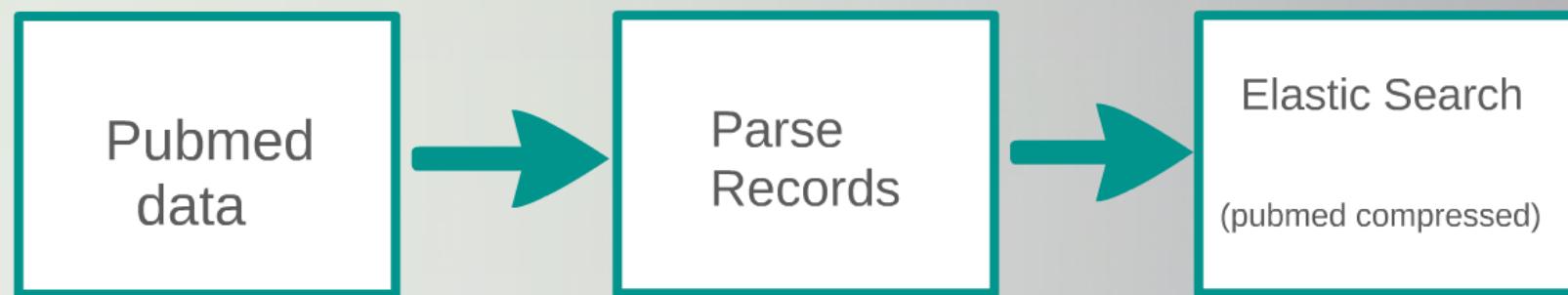
elasticsearch



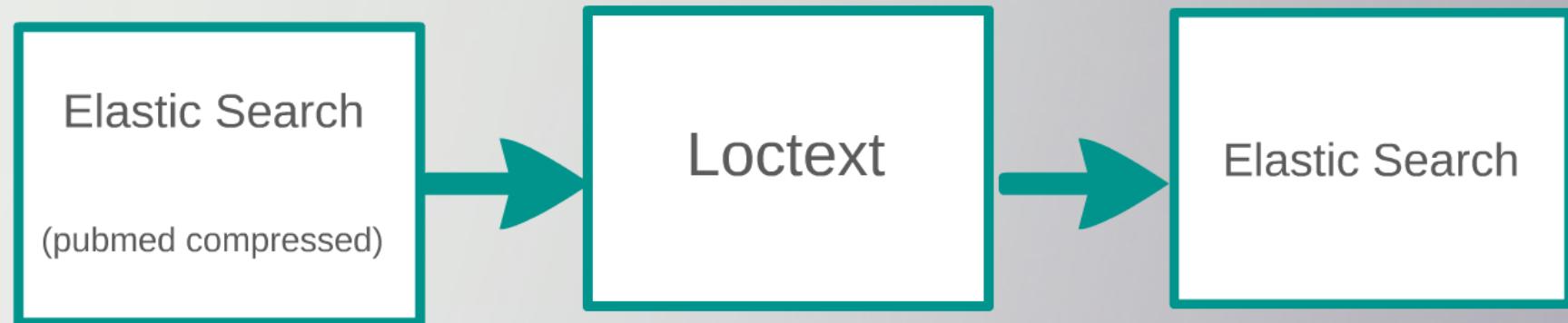
kibana



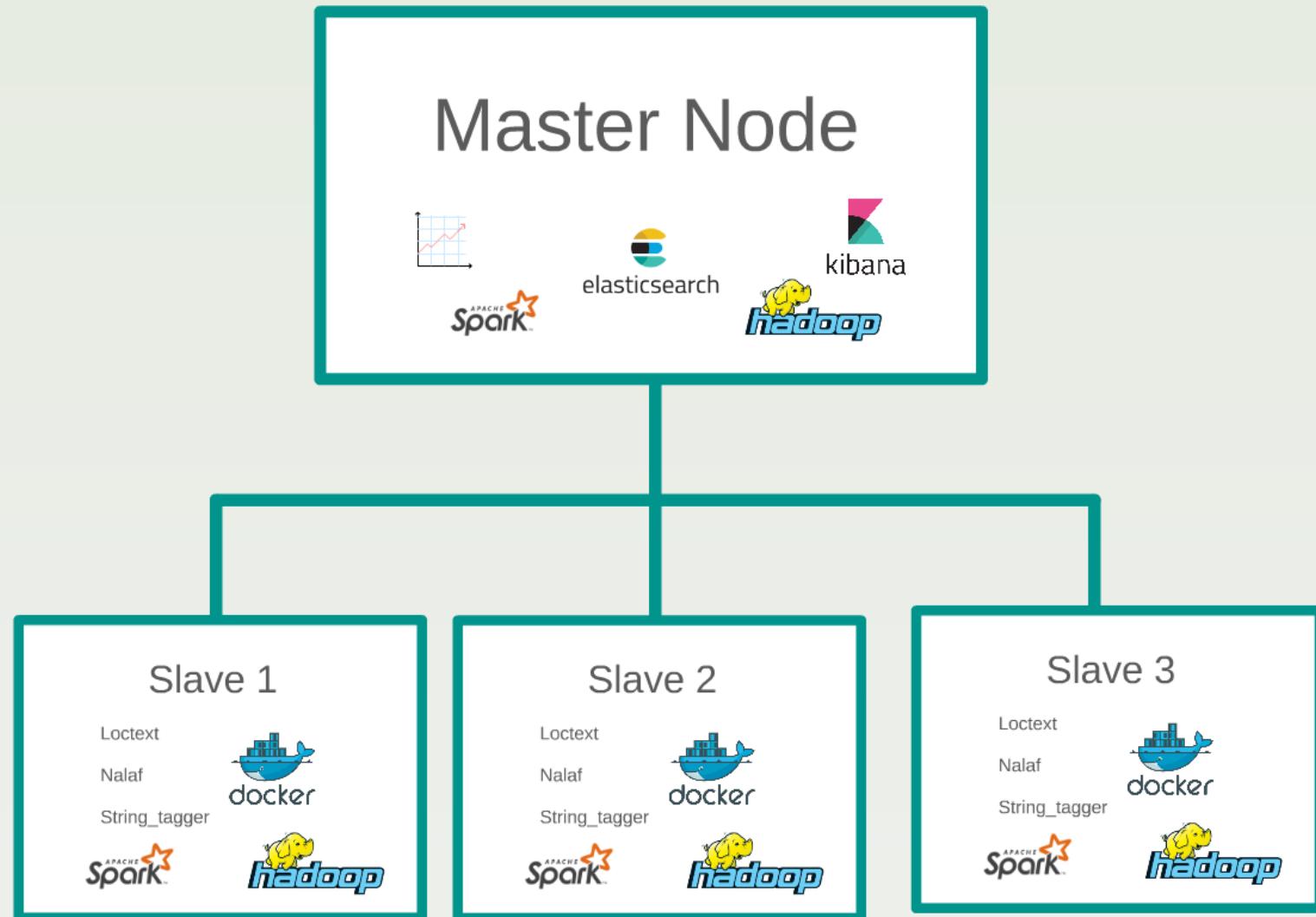
Parsing and storing in Elastic Search



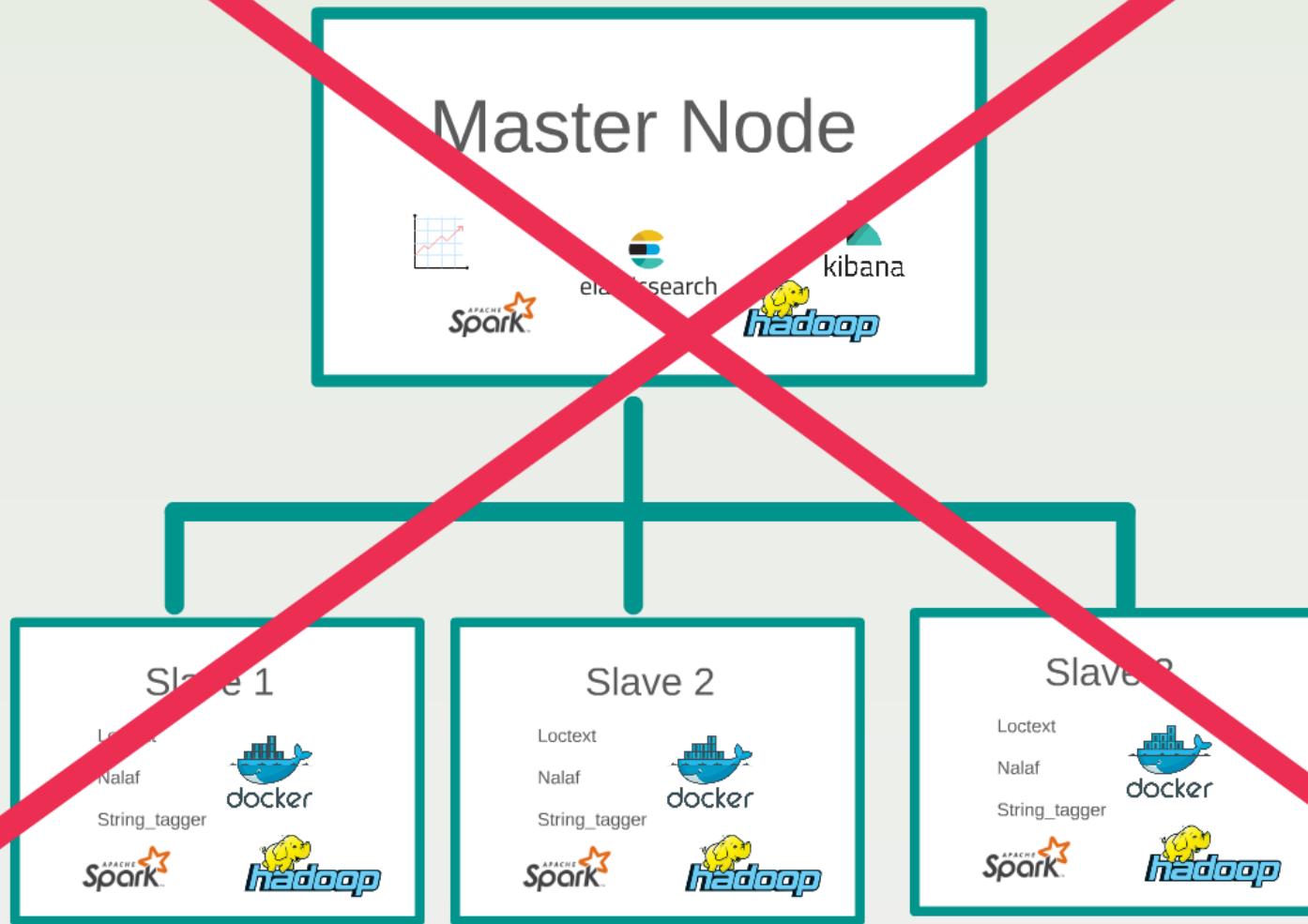
Updating Protein Localization



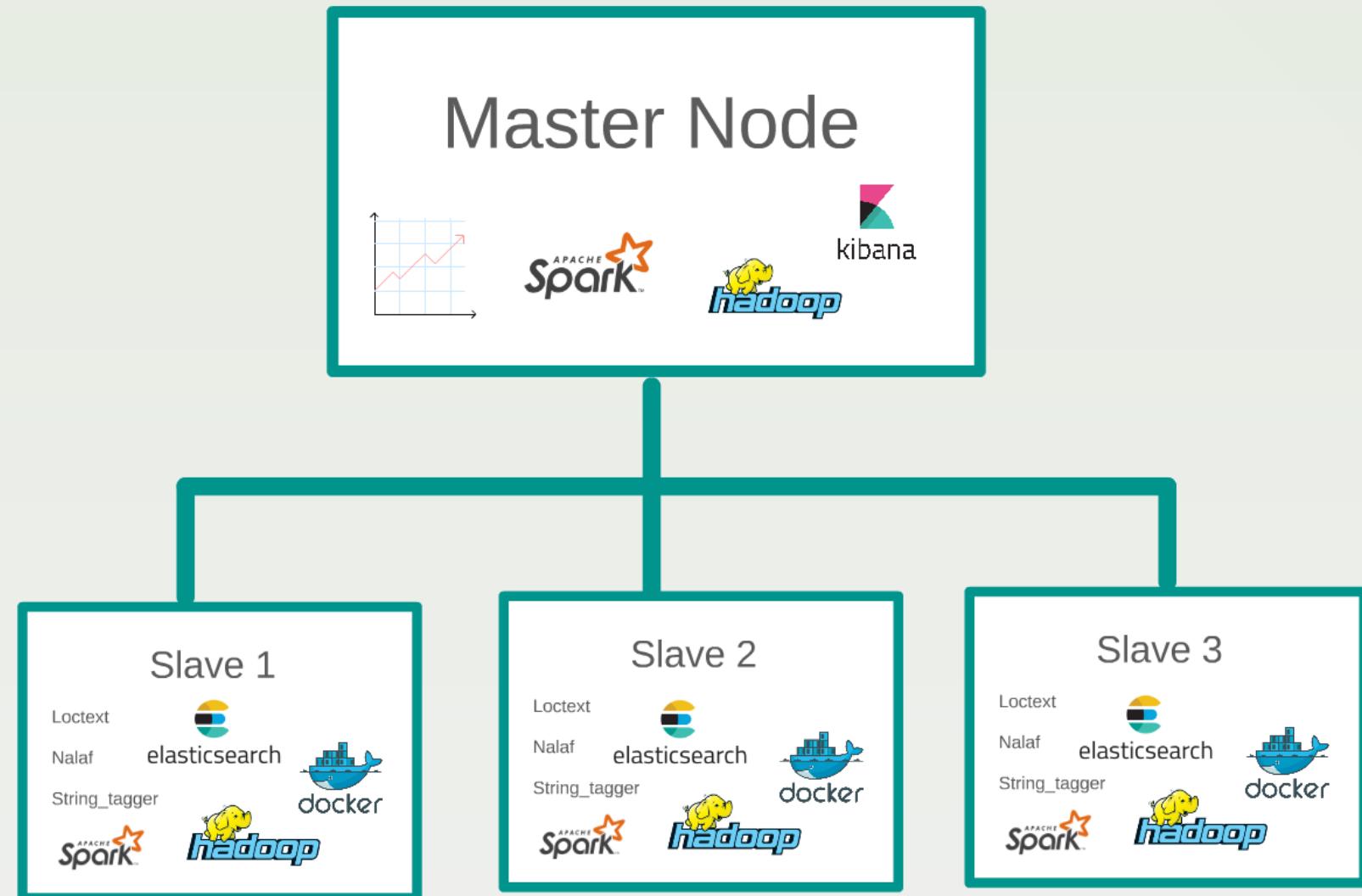
Project Architecture I



Project Architecture I



Project Architecture 2



Visualizing data with



<https://goo.gl/7AdbEz>

Weekly Milestones

Virtual Machines setup



Single and MultiNode Cluster setup



Exploration and programming of Hadoop & Spark



Extraction of NCBI Database and Elsevier Data



LocText

Elastic Search, LocText, Parsing and storage of Full Text Papers



Text-mine the relationship of Protein & Cell with LocText and MapReduce



Obstacles



Port blocking / Live nodes
Java heap space



HTTPs error Efetch API



Docker IO
Docker ip tables error



elasticsearch

Unstructured data structure.
Memory error

LocText

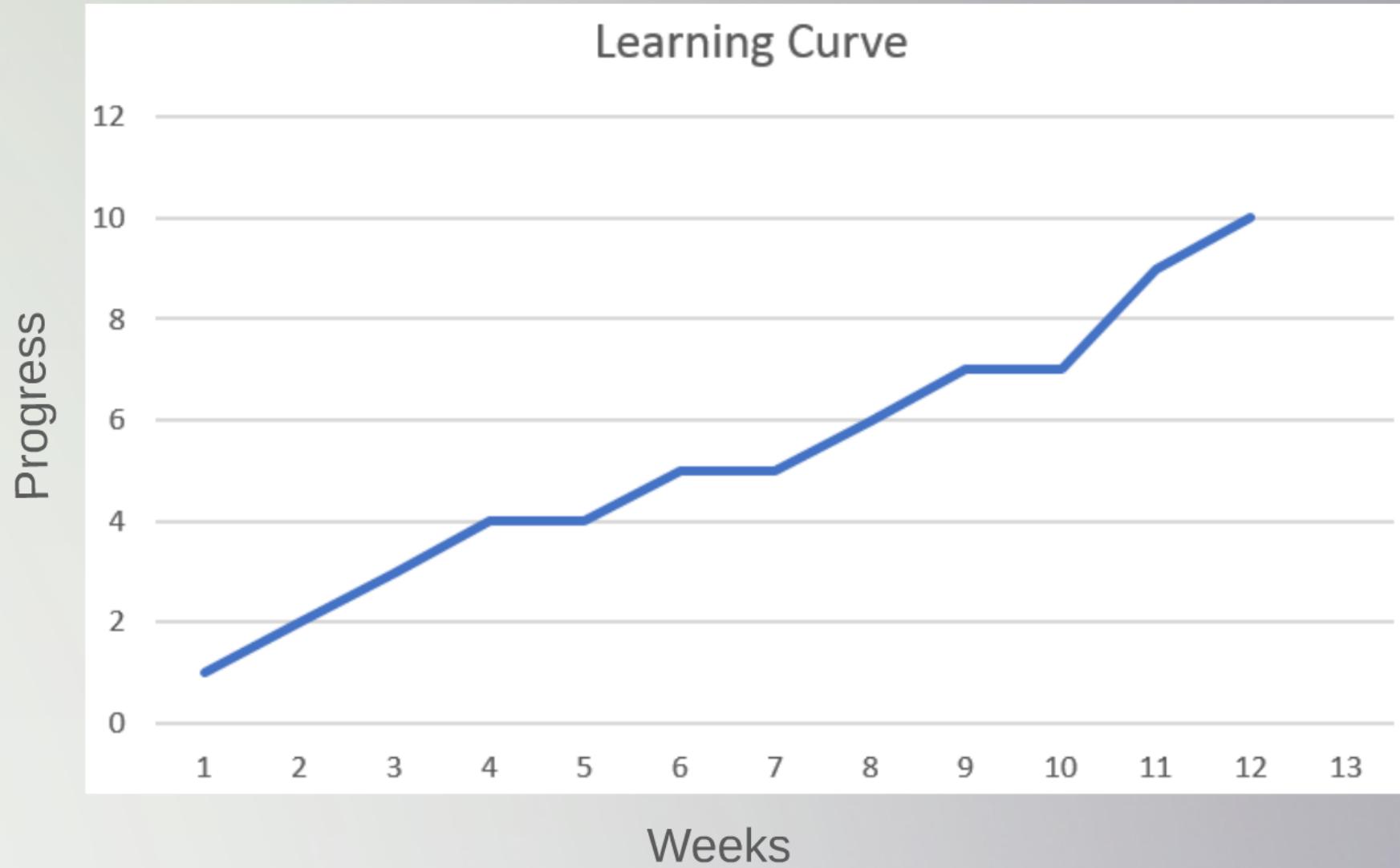
Import Errors (String Tagger)
Input to Loctext is not a protein



Biggest Error

```
@@@@@  
@     WARNING: REMOTE HOST IDENTIFICATION HAS CHANGED!      @  
@@@@@  
IT IS POSSIBLE THAT SOMEONE IS DOING SOMETHING NASTY!  
Someone could be eavesdropping on you right now (man-in-the-middle attack)!  
It is also possible that a host key has just been changed.  
The fingerprint for the ECDSA key sent by the remote host is  
SHA256:MoxygOTJN4TImrj2Mdqyz16XWnpr0WZlSA+gBc2LFGEk.  
Please contact your system administrator.
```

Learning curve



Conclusion

- Introduction to distributed technologies.
- Setting up the cluster from the scratch.
- Distributed programming experience.
- Implemented and Visualized the protein locations.
- Mini talks.

