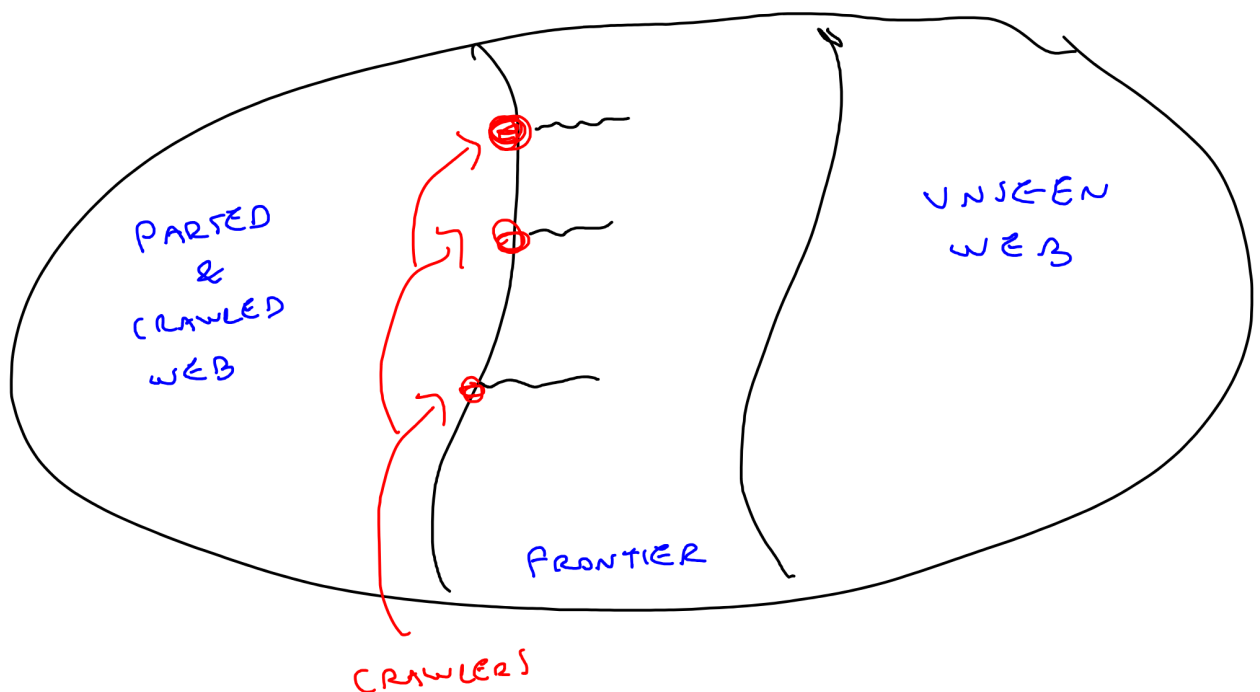


# CRAWLING

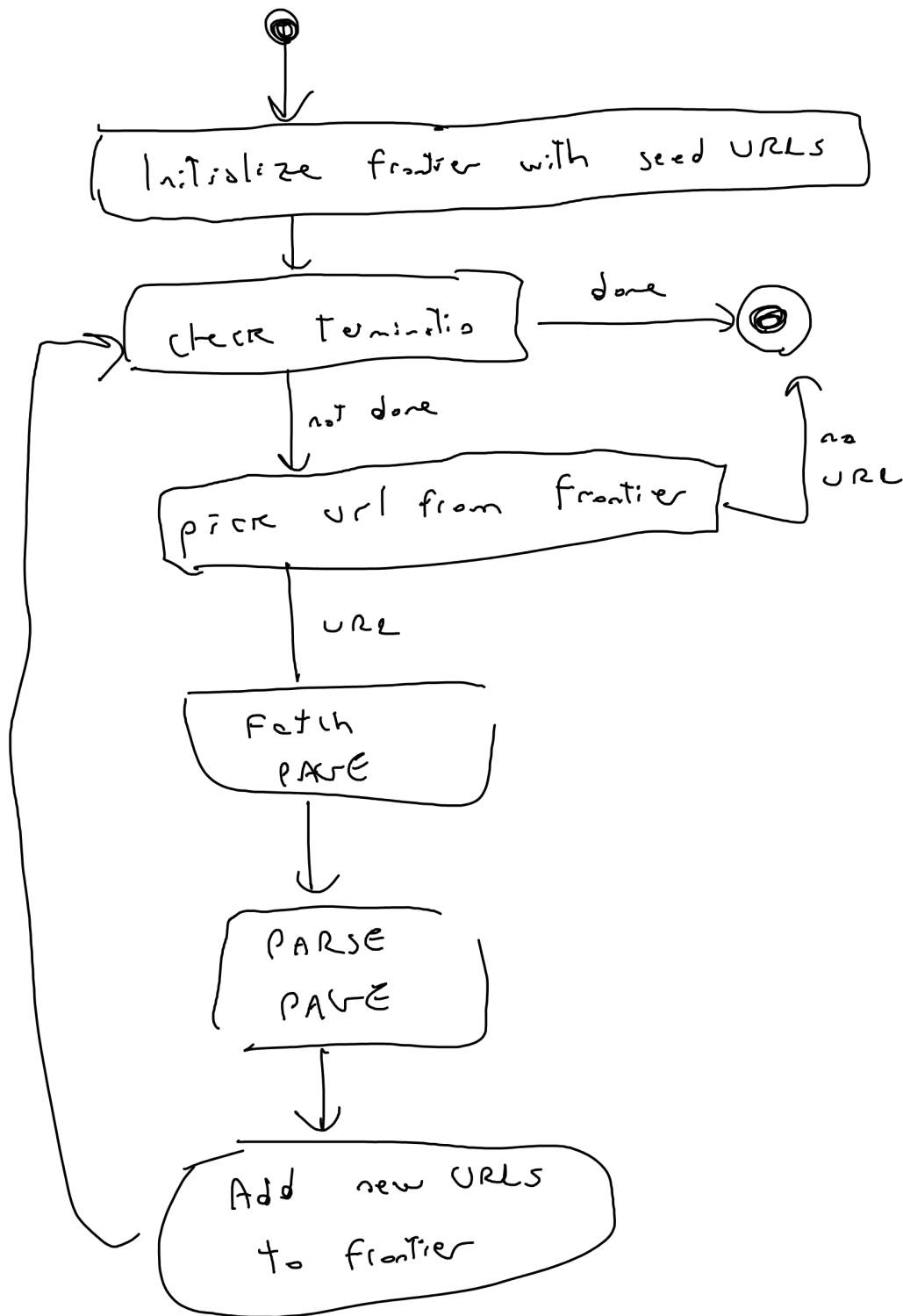
WALKING OVER THE WEB GRAPH IN ORDER  
TO TRACE ITS ROUTES

## CHALLENGES

- You should crawl the best pages first (Quality)
- You should avoid duplication or near duplication (Efficiency)
- Minimize the load of the crawled hosts (Etiquette)
- Avoid spam pages and spider traps (Malicious pages)
- How to cover the whole web (Coverage)
- How much our competitors are able to cover (Relative coverage)
- How often should I crawl (Frequency)



# CRAWLER WORKFLOW



Several metrics to pick up the best url  
from the frontier = Random, BFS, DFS,  
Toprc driven, Random, Page Rank - - -  
↳ how popular is  
the page