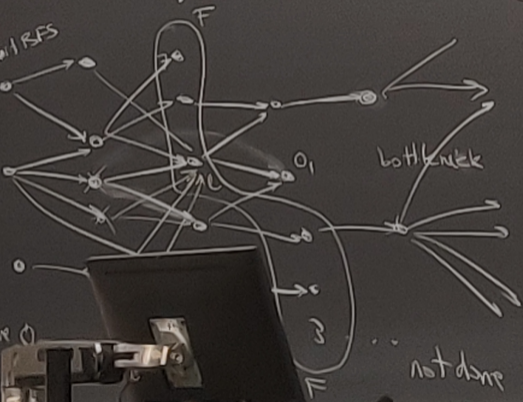
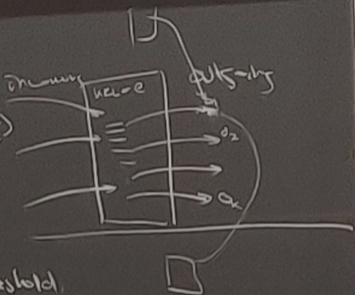


crawler (F = {add URLs})
 queue
 t = frontier queue



- ① URLs = depend (F, batch)
- ② In each LEVEL
 - (A) GET(e) - http head call \Rightarrow meta info
 - wget call \Rightarrow HTML source
 - polite: DELAY 40s, Robots.txt (server)



? CLEAN(e)
 EXTRACT outgoing links (e) $\Rightarrow \{o_1, o_2, \dots, o_k\}$
 ANALYZE CONTENT(e) - relevance (topic) \geq threshold

- diversity / novelty
- near-duplicate (source \approx source' E INDEX)
- page rank \approx importance from links (links)
- domain \approx importance due to site / author.

(B) ? maybe
 ignore (e, source)
 continue FOR (next e)

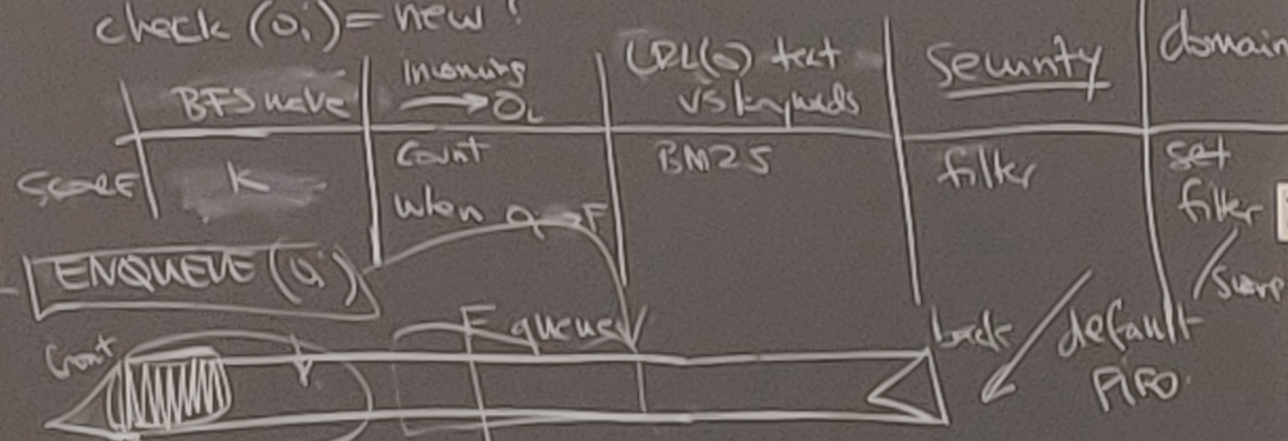
not done

(2C) Store (f, source) \Rightarrow ES

For each $O_i \in \text{outgoing links}(e)$

$O_i = \text{canonical}(v_i)$

check (v_i) = new?



- mechanisms (operations, data structures)

- local resolving/rearrangement intrinsic