

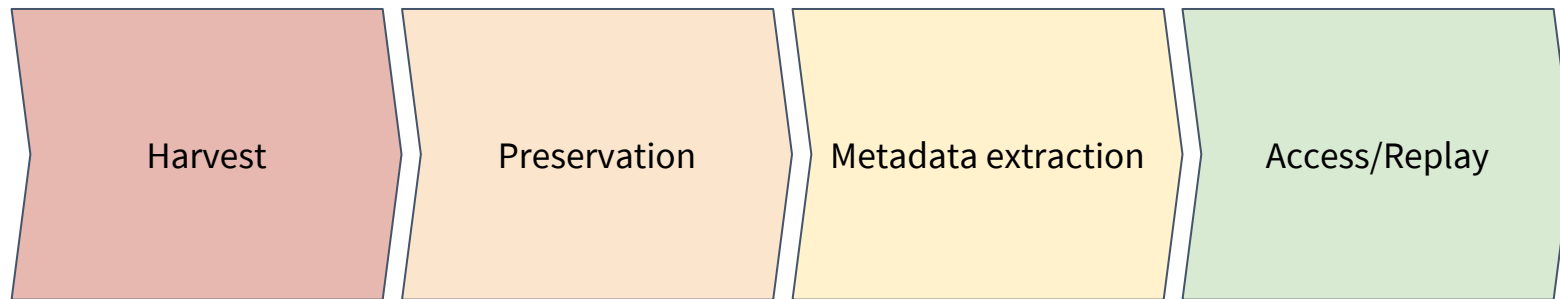


# Sifting Needles out of (Well-Formed) Haystacks: Using LOCKSS Plugins for Web Archive Metadata Extraction

Thib Guicherd-Callin – Technical Manager, LOCKSS Program  
[thib@cs.stanford.edu](mailto:thib@cs.stanford.edu) – [github.com/thibgc](https://github.com/thibgc)

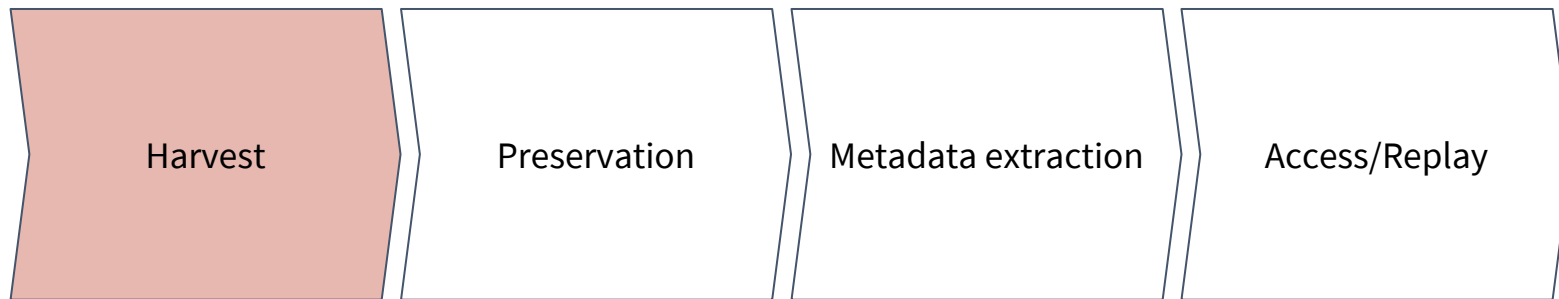


# A Day in the Life





# Harvest

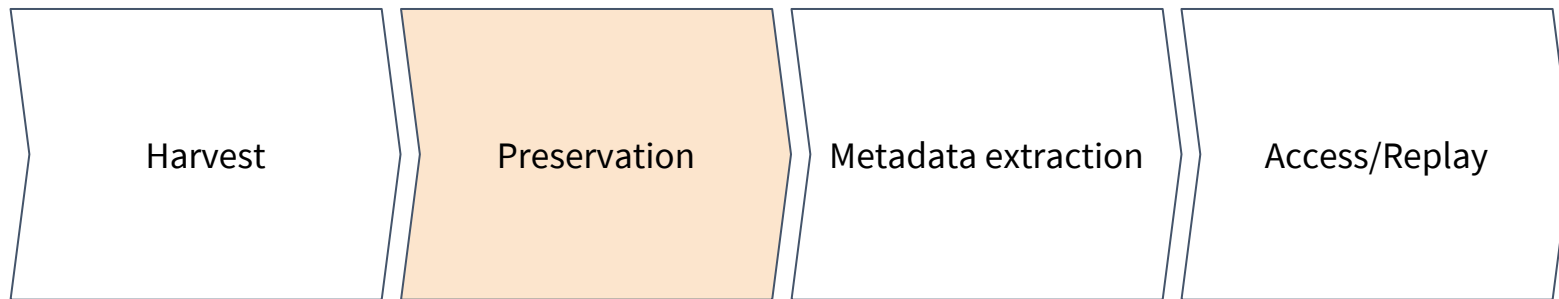


- Permission URLs
- Start URLs
- Crawl rules
- HTTP response handlers
- Content validators
- Crawl filters

- Link extractors
- URL normalizers
- URL consumers
- Permission checkers
- Substance checkers
- Crawl windows
- Fetch interval controls



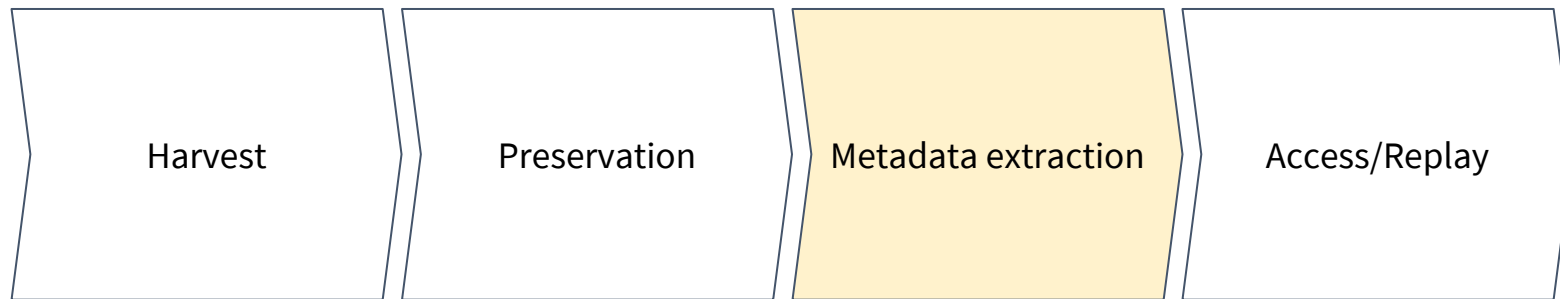
# Preservation



- Hash filters
- URL weighting



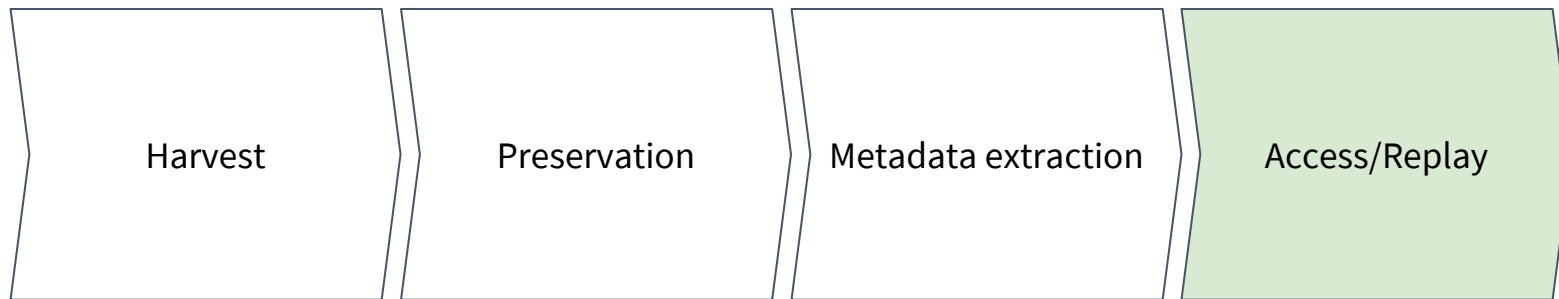
# Metadata Extraction



- Article iterators
- Article metadata extractors
- File metadata extractors



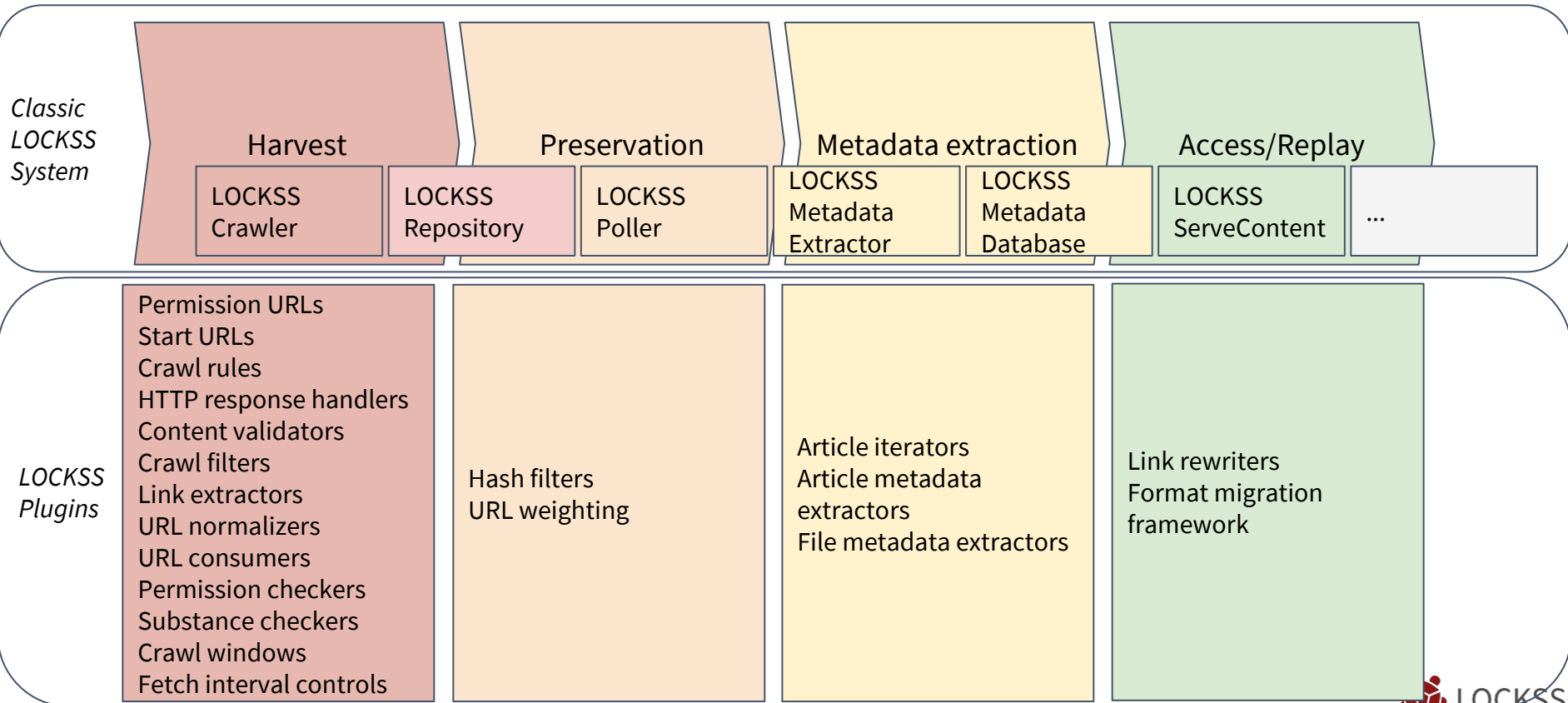
# Access/Replay



- Link rewriters
- Format migration framework

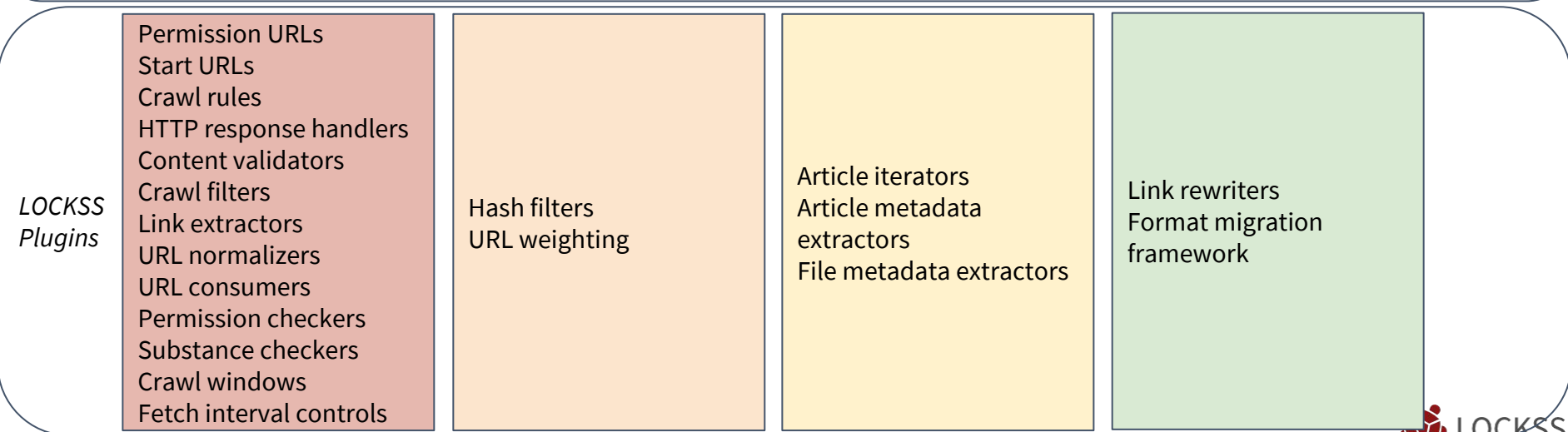
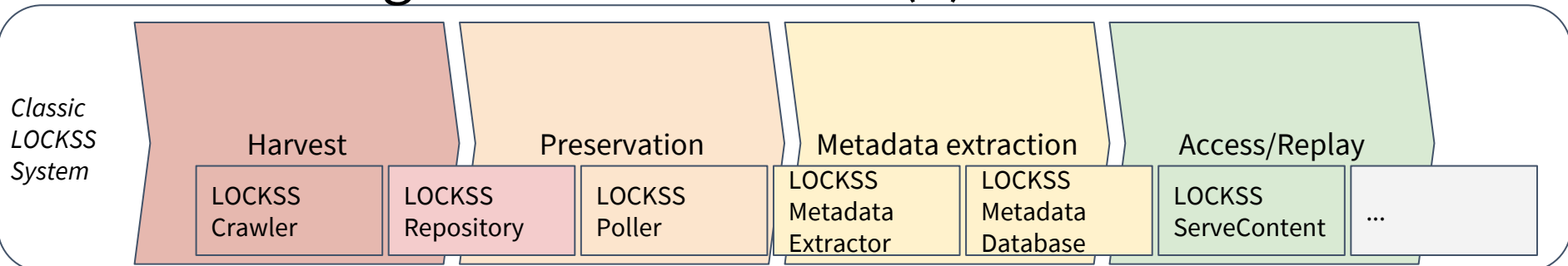


# Classic LOCKSS Architecture





# What's Wrong with This Picture? (1)







# LAAWS Re-Architecture

*Classic  
LOCKSS  
System*

LOCKSS  
Crawler

LOCKSS  
Repository

LOCKSS  
Poller

LOCKSS  
Metadata  
Extractor

LOCKSS  
Metadata  
Database

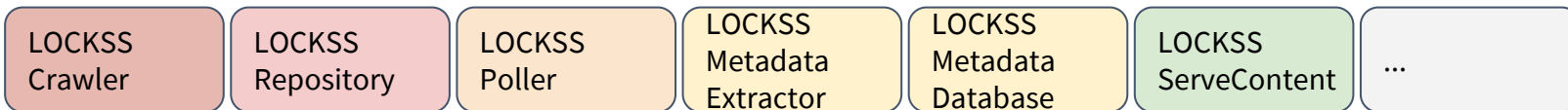
LOCKSS  
ServeContent

...



# LAAWS Re-Architecture

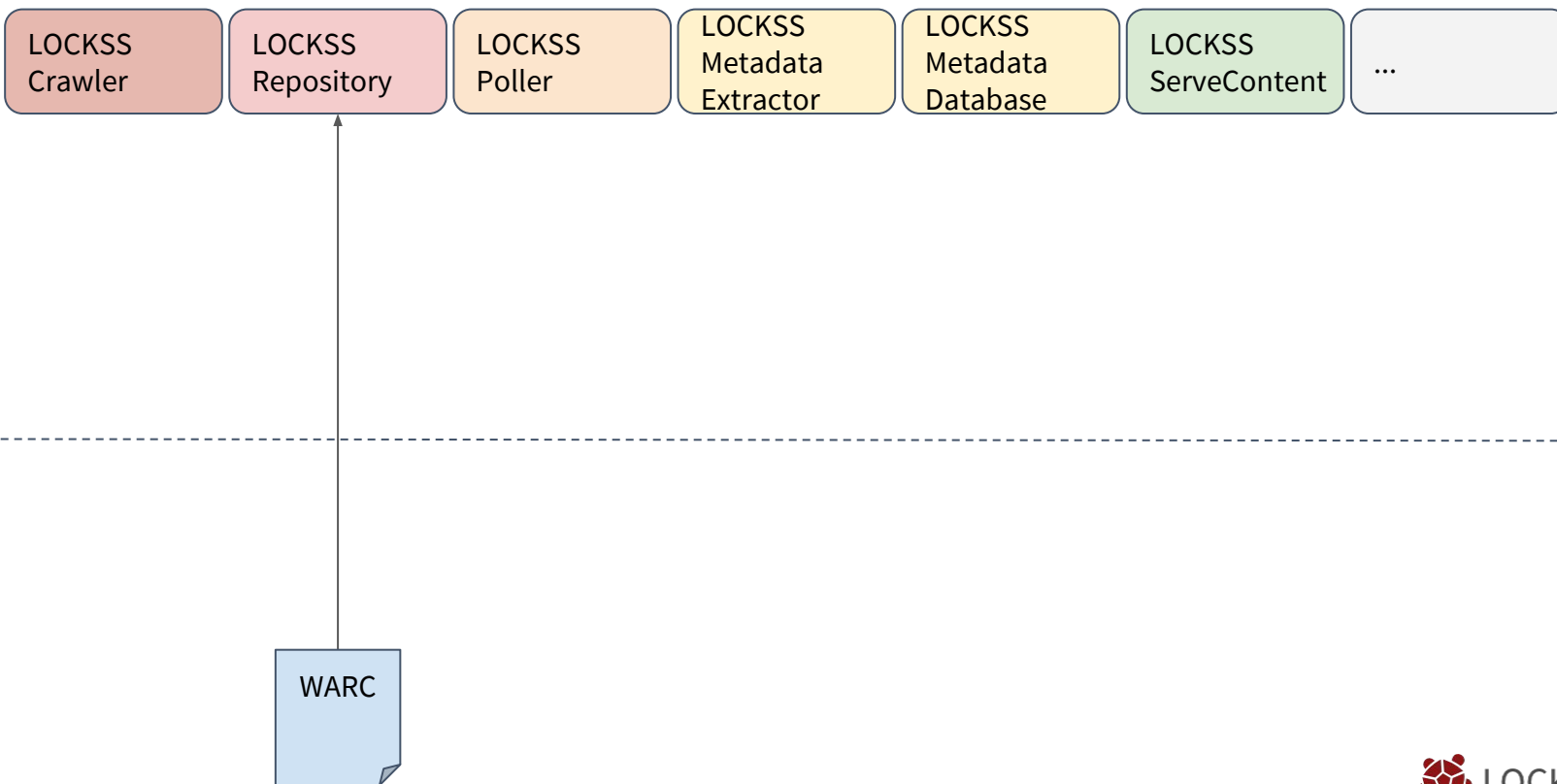
*Re-Architected  
LOCKSS  
System*





# Repository Interoperability

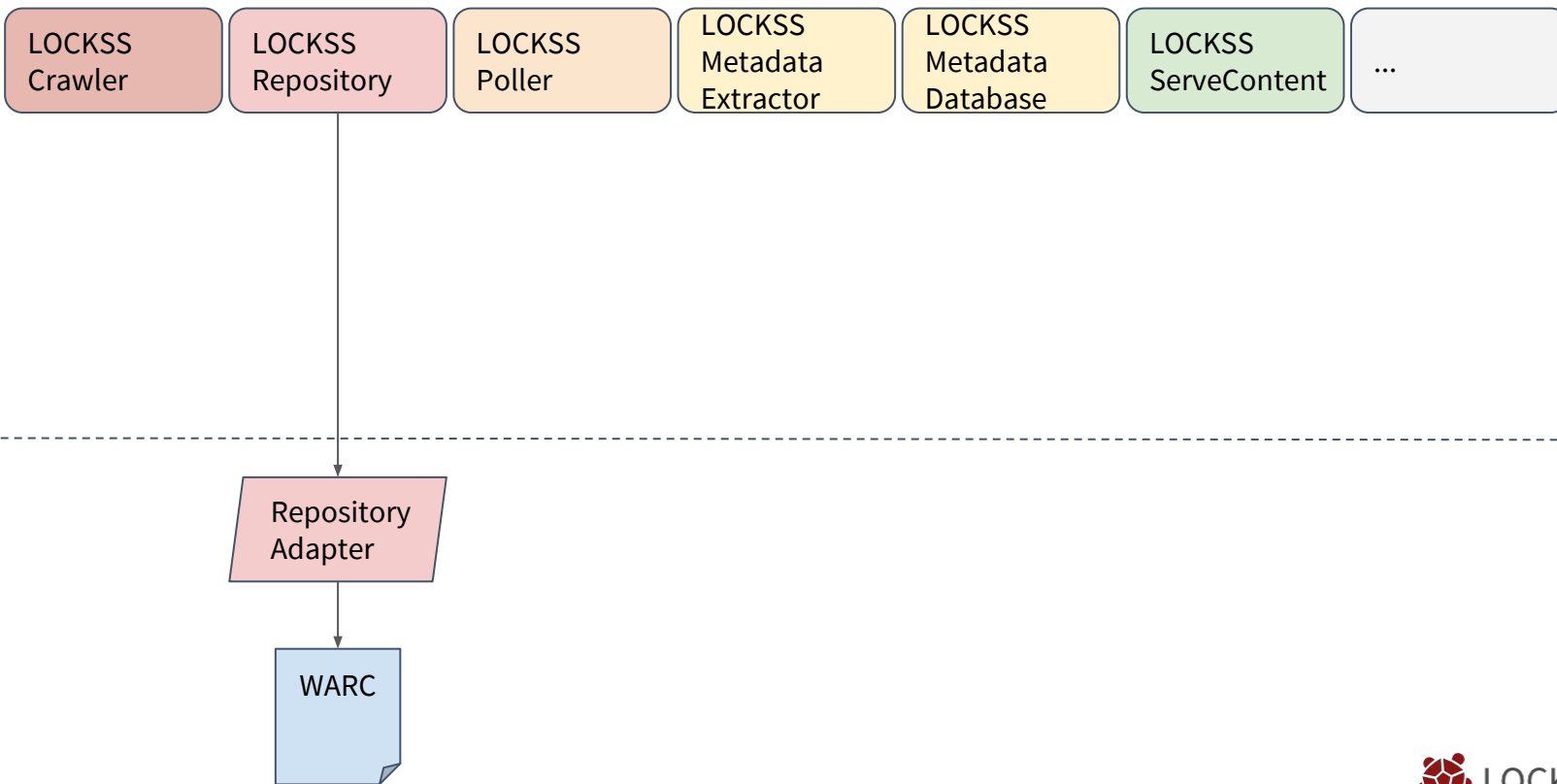
*Re-Architected  
LOCKSS  
System*





# WARC Ingest

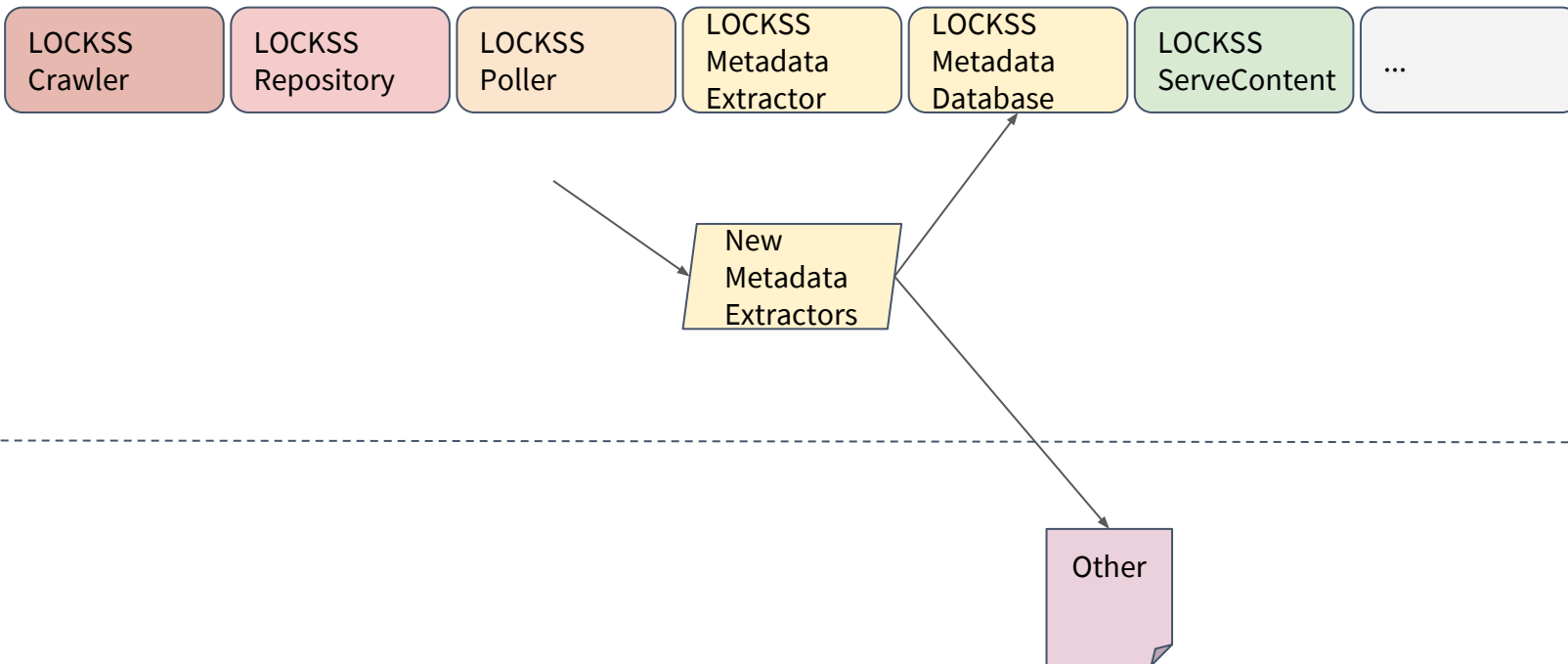
*Re-Architected  
LOCKSS  
System*





# WARC Ingest

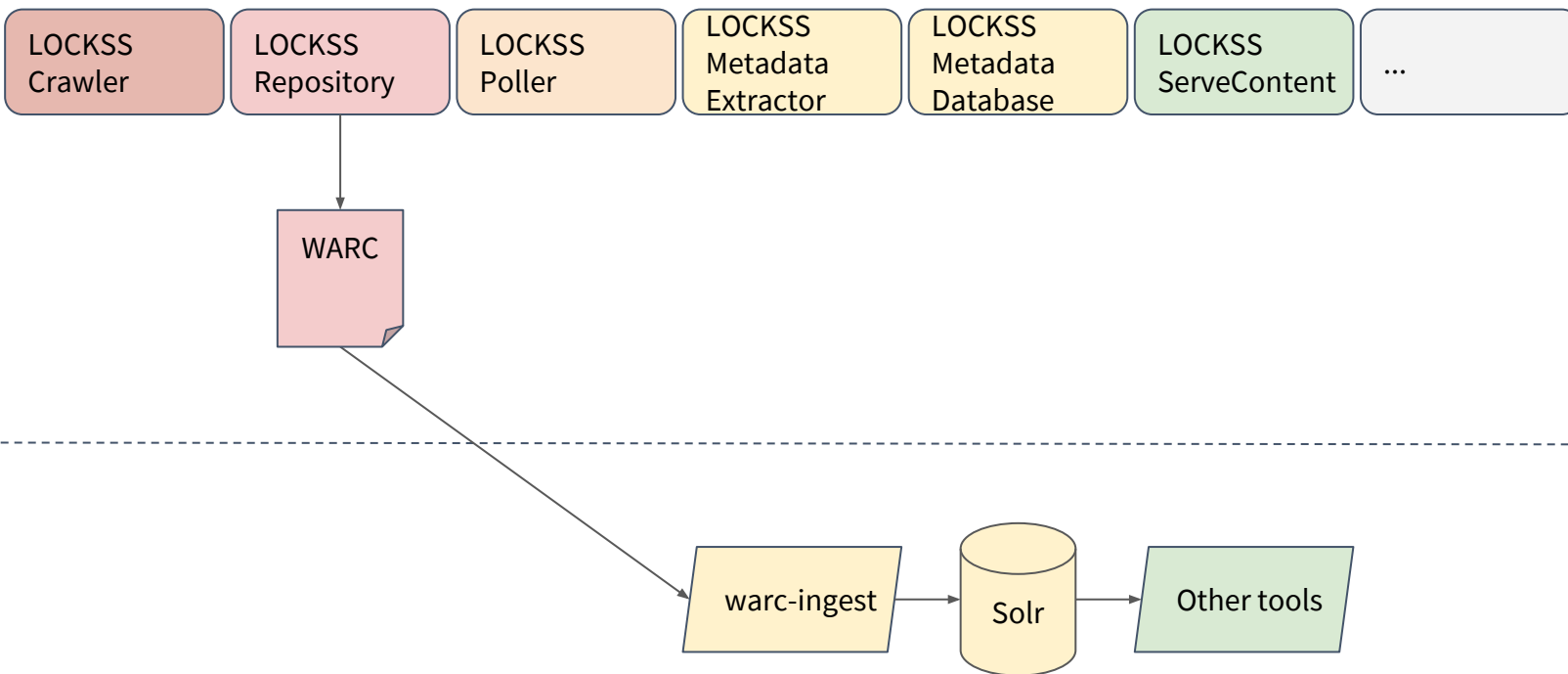
*Re-Architected  
LOCKSS  
System*





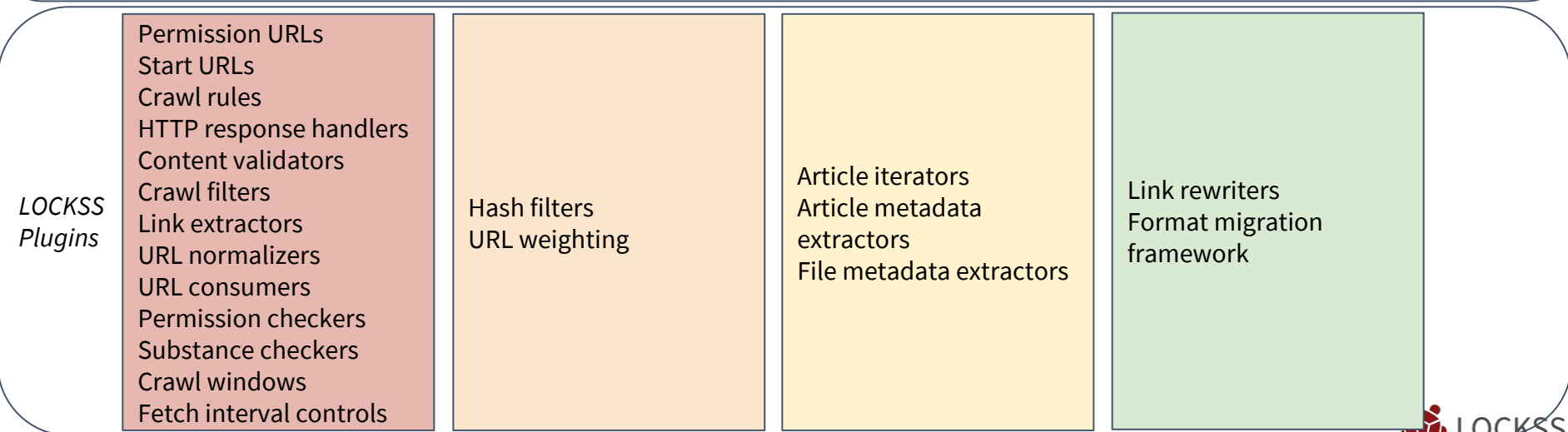
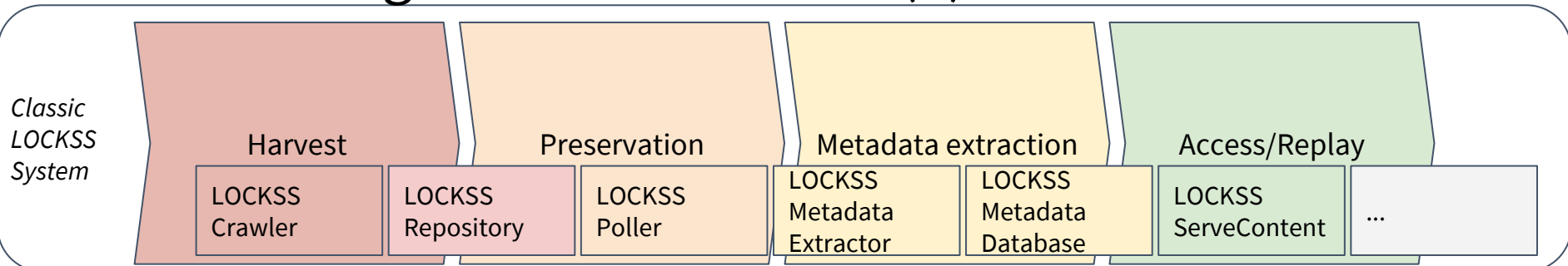
# WARC Ingest

*Re-Architected  
LOCKSS  
System*





# What's Wrong with This Picture? (2)





# Java Tooling from LOCKSS Plugins

- Running metadata extractors externally
- HTML transforms
- PDF transforms
- RIS parser
- Running other things externally:
  - Link extractors
  - Filters
  - Content validators
  - URL normalizers
  - etc.





# Thank You

- Nascent community tools:
  - LOCKSS Documentation Portal: <https://lockss.github.io/>
  - LOCKSS Slack: <https://tinyurl.com/slackjoinlockss>
- Q&A