

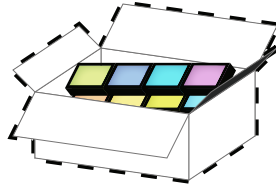
Memory structure: Layers

Layer 4:
User
Corpus

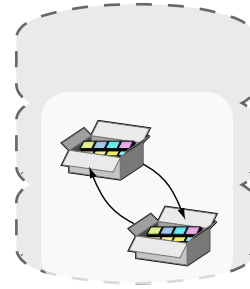


Corpus of WebEntities:
A database of WebEntities (metadata)

Layer 3:
Dynamic
WebEntities
Building

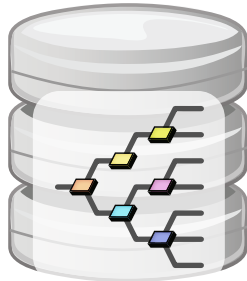


WebEntities:
A WebEntity is one or more LRUs. They are defined in the LRU Index.

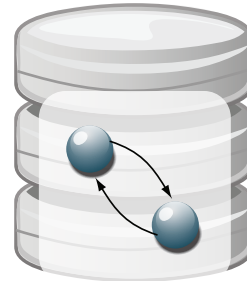


WebEntities Links:
Dynamically build from LRU Index + Nodes Links

Layer 2:
Nodes &
LRU Index



LRU Index:
All LRUs with Node flags and Pages



Nodes Links

Layer 1:
RAW Data



Raw Pages:
Text files that contain the **links** on each page (and possibly more)

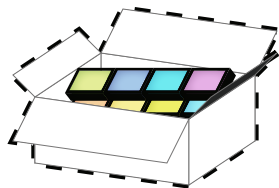
Memory structure: Data Flowing from a Harvest

Layer 4:
User
Corpus

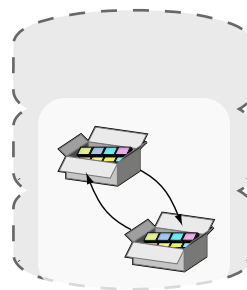


Corpus of
WebEntities

Layer 3:
Dynamic
WebEntities
Building

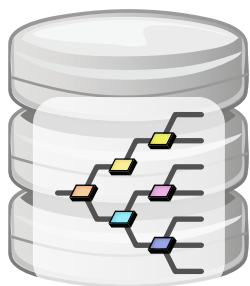


WebEntities



WebEntities
Links

Layer 2:
Nodes &
LRU Index



LRU Index



Nodes Links

3. The LRUs of the Pages, and the corresponding Nodes are stored in the LRU Index.

4. Links between Pages are aggregated as Links between Nodes and stored.

Layer 1:
RAW Data



Raw Pages

1. Harvest: Raw Data comes in and is stored in the Layer 1.

2. The content added to the RAW Data (Pages) is parsed as LRUs and Nodes.

The Web

Memory structure: Data Provided for a Query

User Interface

Layer 4:
User
Corpus

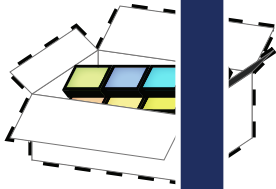


Corpus of
WebEntities

5. Get the metadata
of the neighbor
WebEntities

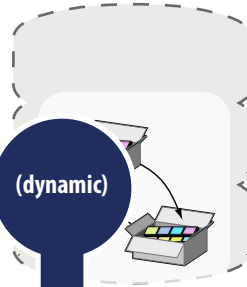
Query: Retrieve
the neighborhood
of a WebEntity

Layer 3:
Dynamic
WebEntities
Building



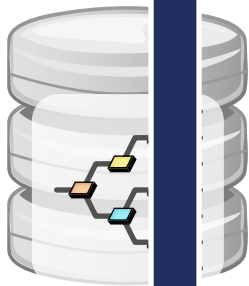
WebEntities

1. Retrieve the
LRUs defining
the WebEntity



WebEntities
Links

Layer 2:
Nodes &
LRU Index



LRU Index

2. Get the
Nodes prefixed
by these LRUs



Nodes Links

3. Get the
neighbor
Nodes

Layer 1:
RAW Data



Raw Pages

4. Get the
WebEntities of the
neighbor Nodes

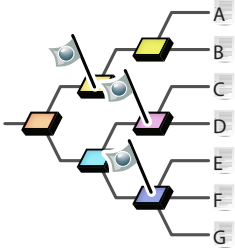
Process: Retrieving the Links of a WebEntity



1. Get WebEntity



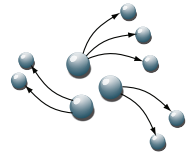
2. Get the one or several LRUs defining it



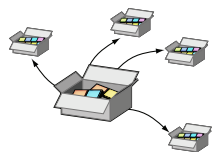
3. From the LRU Index, get all LRUs prefixed by these LRUs (from the WebEntity) that are flagged as Nodes.



4. It defines a certain number of Nodes

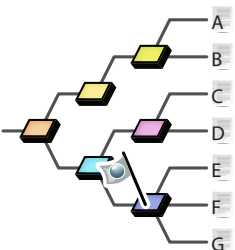


5. Get the Nodes Links from the Layer 2.

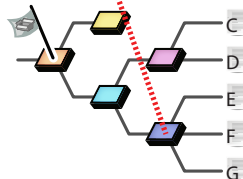


6. From the LRU Index, get the WebEntities of these Nodes (see below).

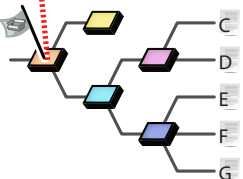
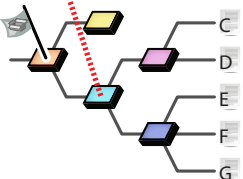
Process: Getting a WebEntity from a Node



1. In the LRU Index, get the LRU of the Node.



2. Search the smallest sub-LRU which is a WebEntity (go back to the root of the tree and get the last WebEntity flag)



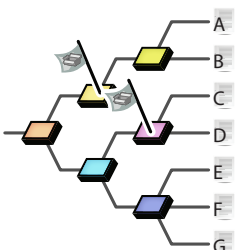
Process: Getting the pages in a WebEntity



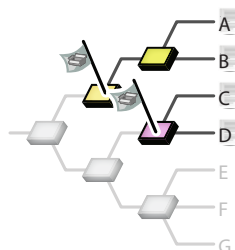
1. Get WebEntity



2. Get the one or several LRUs defining it



3. From the LRU Index, get all LRUs prefixed by these LRUs (from the WebEntity) that are flagged as Pages.



4. In the LRU Index, "walk through the tree" from the LRUs to the last branches that are Pages.