## 1. Extract the chunks using an LM's token probabilities, or using existing human knowledge ...the answer to the ultimate question of life, the universe, and everything is 42 ... ...the humorous tone of his writing is evident, reflecting Douglas Adams' unique wit and creativity the model a better human model itself experts text ...the dry humor are reminiscent of The Hitchhiker's Guide to the Galaxy, reflecting Douglas Adams. corpus **KCD-LM SCD-LM ECD-LM** 3. Inference 2. Build Trie Datastore Search trie -> Match contexts -> Accept or reject chunk -> Generate chunk directly if accepted What is the answer to life the universe and everything? Each Trie has root node as everything reflecting an entry token: everything The answer to everything is reflecting Douglas Adams' 42 Each node represents a text Douglas **chunk** for retrieval: is 42 Each node stores chunk's Adams Adams' preceding contexts: to (everything is 42 reflecting Douglas Adams' The answer