

Start
Data Source
dark lyrics.com
metal music
100

process
web scraper
build_dataset.py
- scrapes A-Z bands
- Extracts meta data

Data
All-lyrics.txt
617,708 lines
~20mb

Decision
Format Choice

option 1

Data
Structured
structured-lyrics.txt
- XML like

Process
Data Processing
data_struct.py
- Cleans raw format
- structures data

Recommended

Data
Training format
Training-format-lyrics.txt
- Best for AI

Not Recommended

Data
Lyrics only
lyrics-only.txt
- too small
for quality output

Data
Selected
Format

process
prepare.py
- tokenize text
- creates vocabulary
- splits train/val

• Train.bin
• Val.bin
• Pkl.bin

Data
Binary data
ready for training

Decision
Training strategy

AI
Test Training
3.38m params

1k iterations

AI
Full training
25.64 params
10k iterations

AI
GPT architecture
model.py
- Transforms layers
- self attention
- character level

