



CompositeDatasetDriver

CompositeDatasetDriver

CompositeDatasetDriver

collect file names

collect file names

collect file names

shuffle (seed)

shuffle (seed)

shuffle (seed)

split by rank

split by rank

split by rank

read file contents

read file contents

read file contents

shuffle (seed)

shuffle (seed)

shuffle (seed)

split by worker

split by worker

split by worker

normalize data
dicts

normalize data
dicts

normalize data
dicts

multiplexing

Tokenize

Pad or stuff
content

log meta
data

filter keys

batch

to tensor

to torch
iterable

torch data
loader