

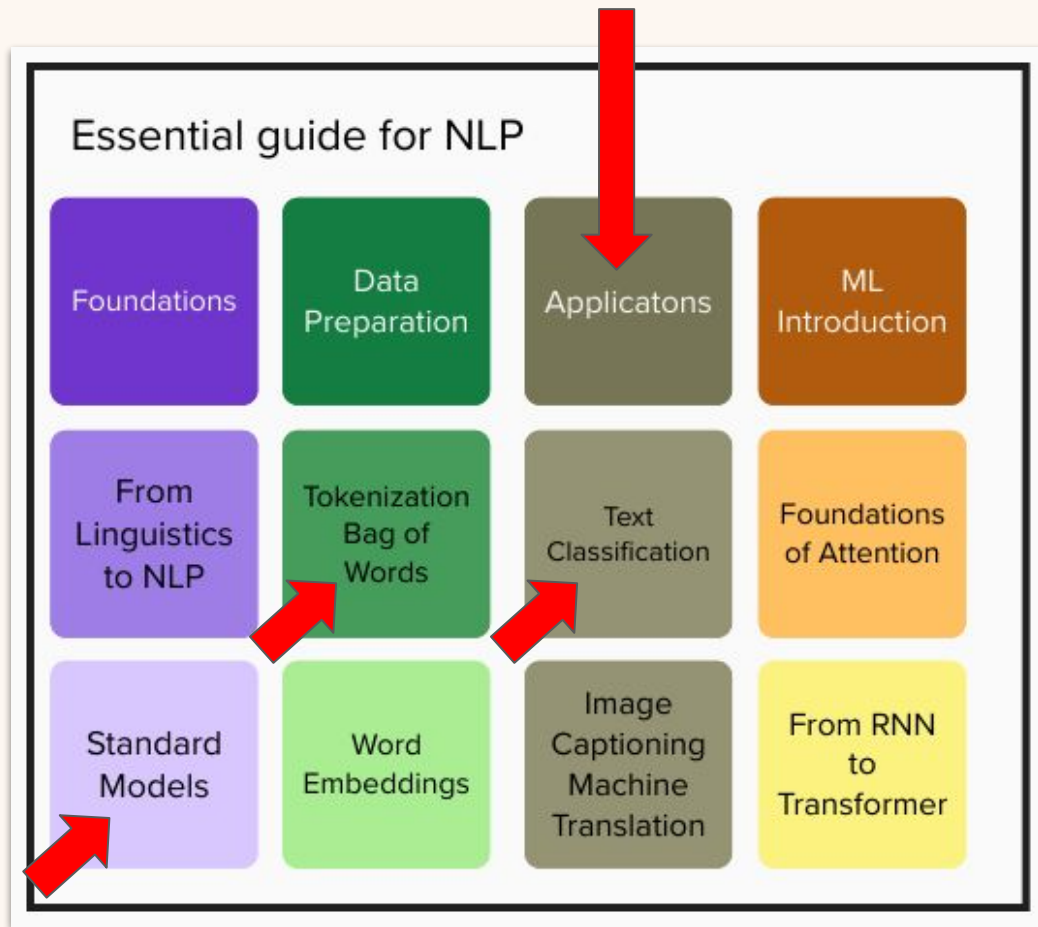
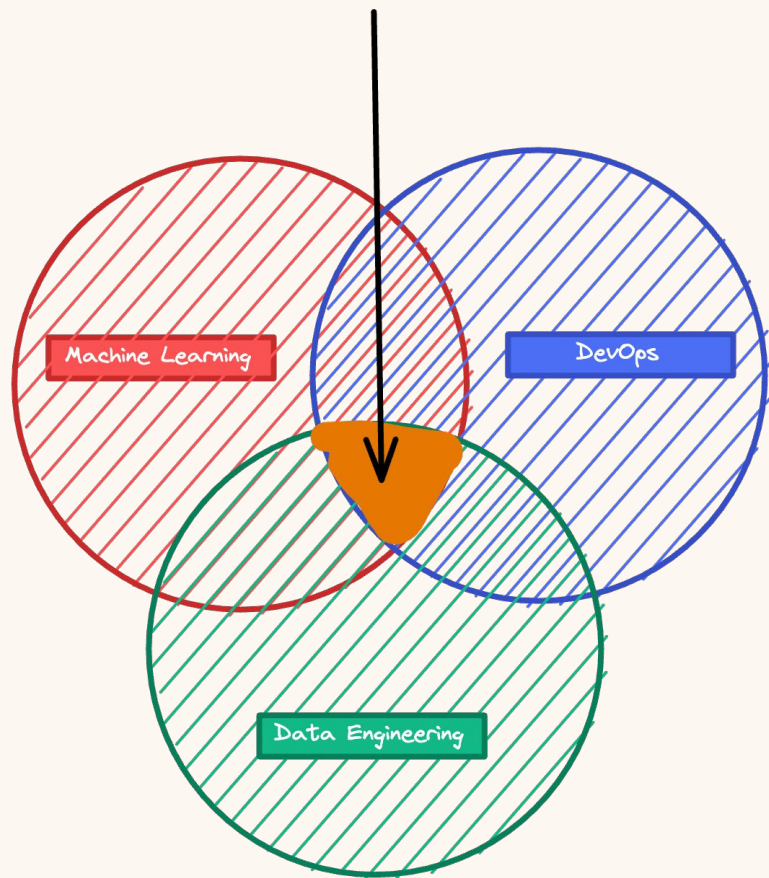
# Essential Guide for NLP

Steps to Process Film Review Data for  
Sentiment Analysis







DCA0305  
ivanovitch.silva@ufrn.br

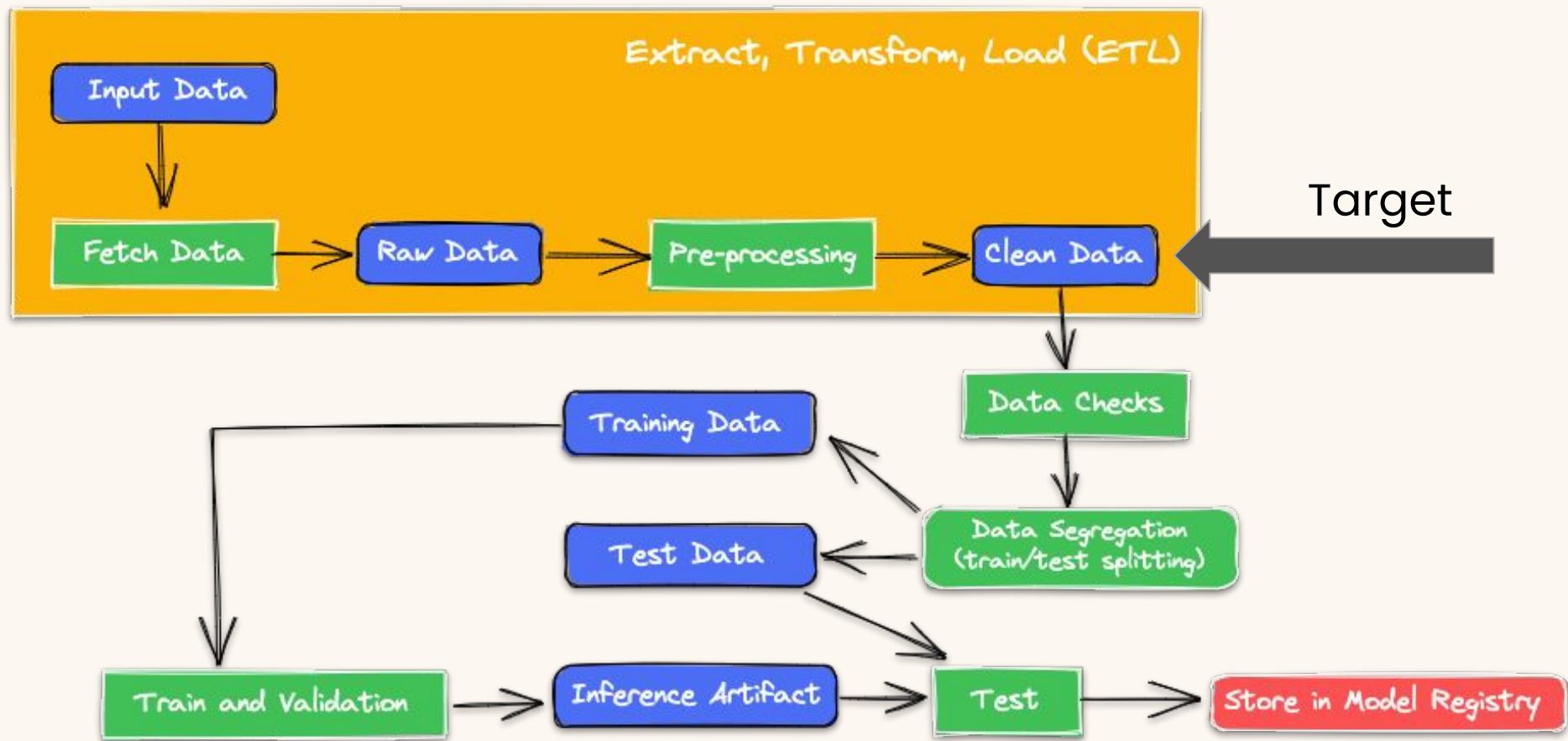


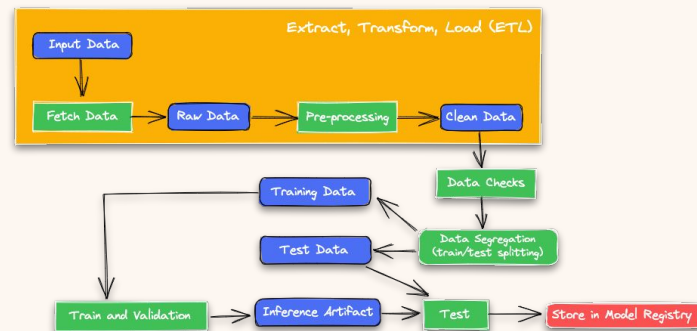
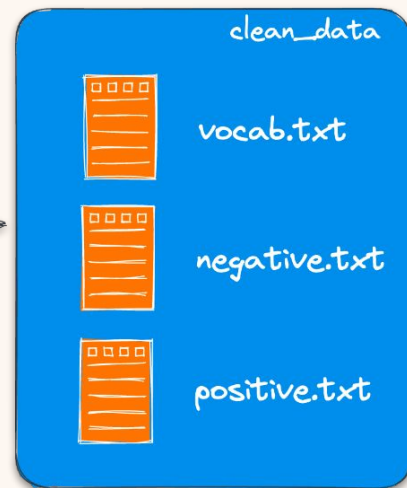
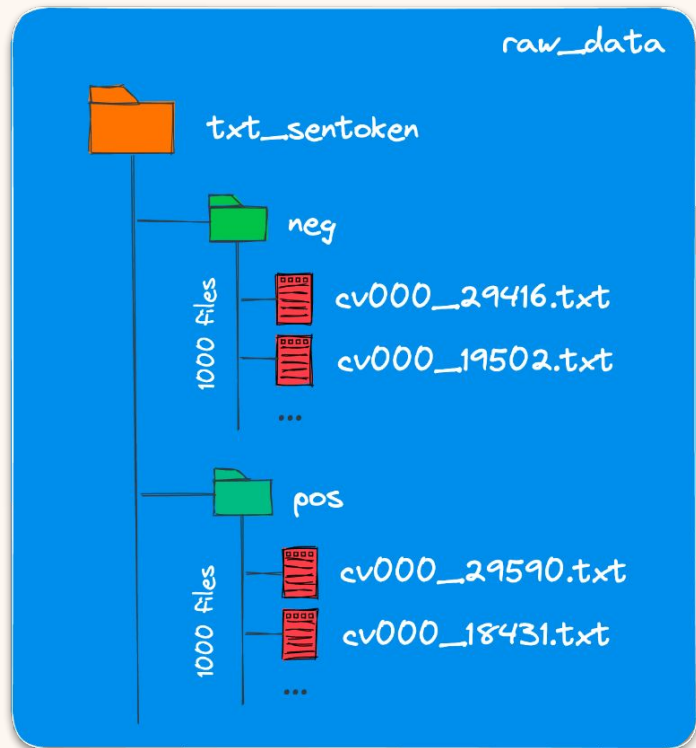
# MLOps



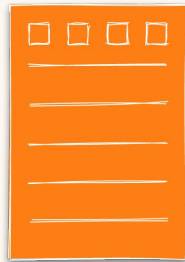
## Top 10 on IMDb this week

 <p>FROM THE CREATORS OF THE HAUNTING OF HILL HOUSE</p> <p>NEVERMORE THE FALL OF THE HOUSE OF USHER</p> <p>ONLY ON NETFLIX   OCTOBER 12</p> <p>★ 8.1 ☆</p> <p>1. The Fall of the House of Usher</p> <p>+ Watchlist</p> <p>▶ Trailer</p>	 <p>FROM THE WORLD OF BOYFRIENDS</p> <p>GEN V</p> <p>prime   SEPTEMBER 29 New Series</p> <p>★ 8.0 ☆</p> <p>2. Gen V</p> <p>Watch options</p> <p>▶ Trailer</p>	 <p>LOKI SEASON 2</p> <p>OCTOBER 6 Disney+</p> <p>★ 8.2 ☆</p> <p>3. Loki</p> <p>+ Watchlist</p> <p>▶ Trailer</p>	 <p>REPTILE</p> <p>JUSTIN TIMBERLAKE ALCIA SILVERSTONE</p> <p>SHED THE LIES</p> <p>ONLY ON NETFLIX   SEPTEMBER 29</p> <p>★ 6.8 ☆</p> <p>4. Reptile</p> <p>+ Watchlist</p> <p>▶ Trailer</p>	 <p>FAIR PLAY</p> <p>PHOEBE DYNOR ALEXANDER ROBERTSON</p> <p>IN SELECT THEATERS SEPTEMBER 13</p> <p>★ 6.5 ☆</p> <p>5. Fair Play</p> <p>+ Watchlist</p> <p>▶ Trailer</p>	 <p>TOTALLY KILLER</p> <p>MURDER IS SO FUN</p> <p>prime   OCTOBER 6 New Movie</p> <p>★ 6.6 ☆</p> <p>6. Totally Killer</p> <p>Watch options</p> <p>▶ Trailer</p>
---	--	---	--	--	--



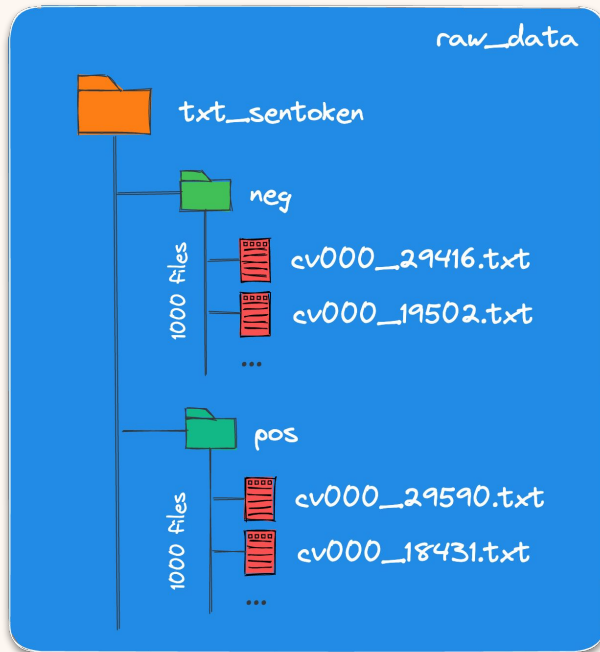




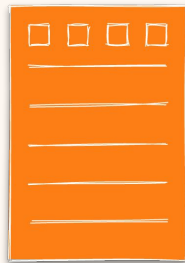


vocab.txt

01 synopsis  
02 meliassa  
03 woman  
04 likes  
05 smoke  
...  
...  
14801 sade  
14802 mongkut  
14802 rumpo

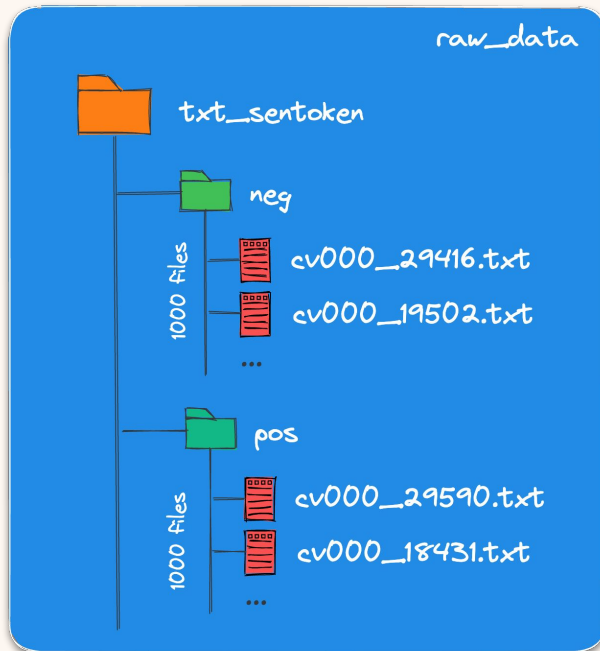


1. Opening and Reading Files
2. Cleaning the Content
3. Compiling a Preliminary Vocabulary
4. Processing Multiple Files
5. Refining the Vocabulary
6. Saving the Final Vocabulary



vocab.txt

01 synopsis  
02 meliassa  
03 woman  
04 likes  
05 smoke  
...  
...  
14801 sade  
14802 mongkut  
14802 rumpo



1. Opening and Reading Files
2. Cleaning the Content
3. Compiling a Preliminary Vocabulary
4. Processing Multiple Files
5. Refining the Vocabulary
6. Saving the Final Vocabulary

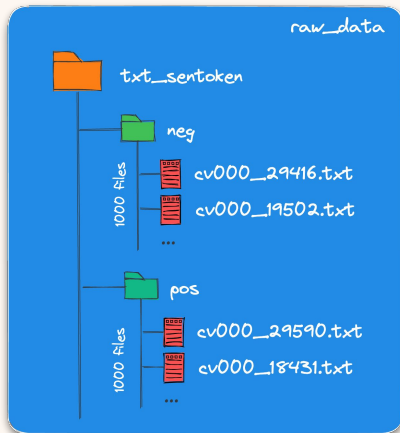
```
# load doc into memory
def load_doc(filename):

    # open the file as read only
    file = open(filename, 'r')
    # read all text
    text = file.read()
    # close the file
    file.close()
    return text
```



vocab.txt

01 synopsis  
02 meliassa  
03 woman  
04 likes  
05 smoke  
...  
14801 sade  
14802 mongkut  
14802 rumpo



1. Opening and Reading Files
2. Cleaning the Content
3. Compiling a Preliminary Vocabulary
4. Processing Multiple Files
5. Refining the Vocabulary
6. Saving the Final Vocabulary

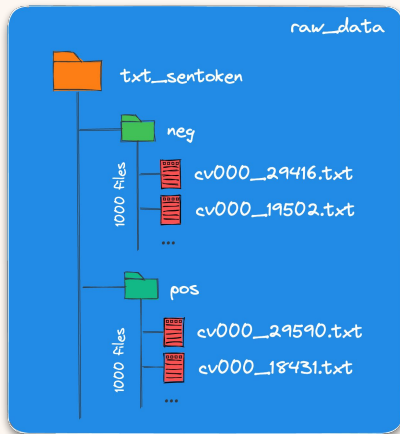
```
# turn a doc into clean tokens
def clean_doc(doc):
    # split into tokens by white space
    tokens = doc.split()
    # prepare regex for char filtering
    re_punc = re.compile('%s' % re.escape(string.punctuation))
    # remove punctuation from each word
    tokens = [re_punc.sub('', w) for w in tokens]
    # remove remaining tokens that are not alphabetic
    tokens = [word for word in tokens if word.isalpha()]
    # filter out stop words
    stop_words = set(stopwords.words('english'))
    tokens = [w for w in tokens if not w in stop_words]
    # filter out short tokens
    tokens = [word for word in tokens if len(word) > 1]
    return tokens
```





vocab.txt

01 synopsis  
02 meliassa  
03 woman  
04 likes  
05 smoke  
...  
14801 sade  
14802 mongkut  
14802 rumpo



1. Opening and Reading Files
2. Cleaning the Content
3. Compiling a Preliminary Vocabulary
4. Processing Multiple Files
5. Refining the Vocabulary
6. Saving the Final Vocabulary

```
# load doc and add to vocab
def add_doc_to_vocab(filename, vocab):
    # load doc
    doc = load_doc(filename)
    # clean doc
    tokens = clean_doc(doc)
    # update counts
    vocab.update(tokens)
```

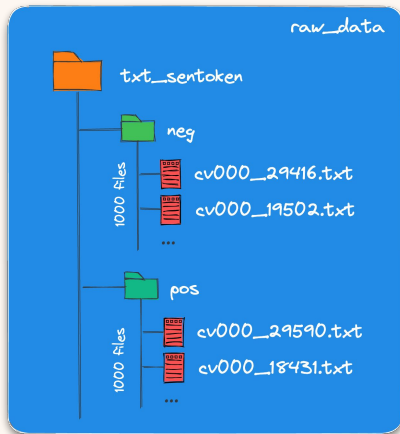
```
# load all docs in a directory
def process_docs(directory, vocab):
    # walk through all files in the folder
    for filename in listdir(directory):
        # skip files that do not have the right extension
        if not filename.endswith(".txt"):
            continue
        # create the full path of the file to open
        path = directory + '/' + filename
        # add doc to vocab
        add_doc_to_vocab(path, vocab)
```

```
# define vocab
vocab = Counter()
# add all docs to vocab
process_docs('txt_sentoken/neg', vocab)
process_docs('txt_sentoken/pos', vocab)
```



vocab.txt

01 synopsis  
02 meliassa  
03 woman  
04 likes  
05 smoke  
...  
14801 sade  
14802 mongkut  
14802 rumpo

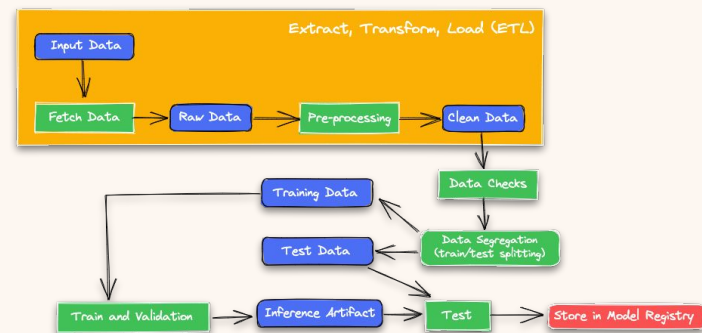
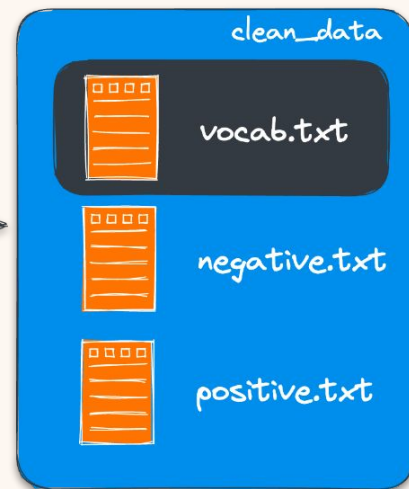
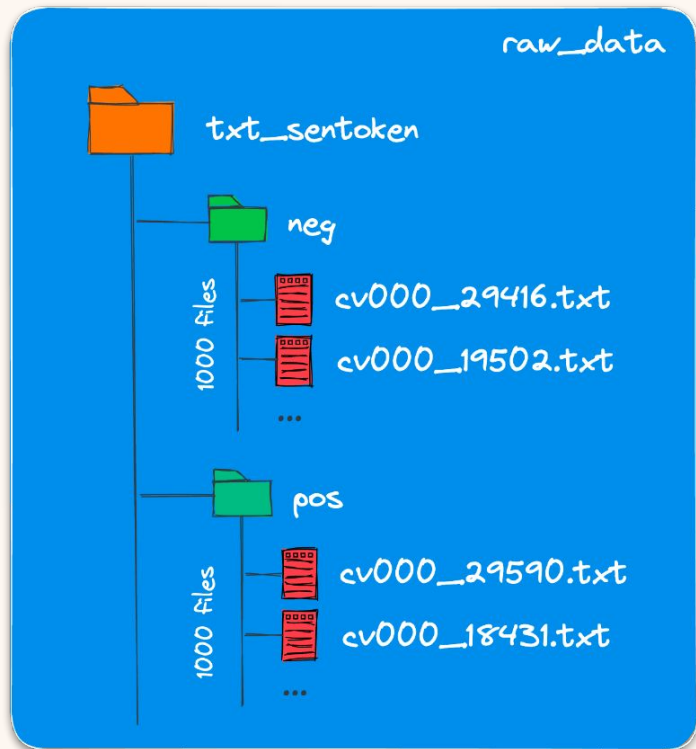


1. Opening and Reading Files
2. Cleaning the Content
3. Compiling a Preliminary Vocabulary
4. Processing Multiple Files
5. Refining the Vocabulary
6. Saving the Final Vocabulary

```
# save list to file
def save_list(lines, filename):
    data = '\n'.join(lines)
    file = open(filename, 'w')
    file.write(data)
    file.close()
```

```
# keep tokens with > 5 occurrence
min_occurrence = 5
tokens = [k for k,c in vocab.items() if c >= min_occurrence]
print(len(tokens))
```

```
# save tokens to a vocabulary file
save_list(tokens, 'vocab.txt')
```



1 already feel hate letters pouring one folks loved wedding singer gets worse much mention titanic sentence he  
2 psychic wounds family sitcom opening segment something film theres guys voice telling us tries imagine biolo  
3 robert redford good playing characters incredible gifts able act like ordinary people natural played fallen  
4 fully loaded entertainment review website coming july didnt really expect much rented stuart saves family mo  
5 fugitive probably one greatest thrillers ever made takes realistic believable characters tells exciting stor  
6 tired hot new releases gone time get video store com dedicated finding hidden gems lie shelves com reviews i  
7 first thing notice movie cold placed minnesota north winter many scenes take place outside long scenes snowc  
8 may years steven spielbergs success jaws years francis coppolas godfather risky ambitious young director nam  
9 think first thing reviewer mention fan xfiles first let assure prior experience series required fully enjoy  
10 present day three sisters reflect parents relationship trying define writerdirector tran scent green vertica  
11 perhaps time say little reading habits really like read ive enjoyed many books lifetime problem im slow read  
12 albert brooks saves day nick time poor summer movies brooks audiences looking cheer positive way may finest  
13 clue unfairly ignored comedy similar murder death big screen version classic board game whats next motion pi  
14 people enjoy science fiction often faced unpleasant surprises due novels stories comic books movies often sc  
15 lisa high art intelligent quiet drama strongest quality aside topnotch central performances perceptive way f

990 fortunate enough attend advance screening upcoming thriller conspiracy theory course big deal reviewing movi  
991 man presented us henry portrait serial killer comes wild tale set within elite white trash south coast plot  
992 steven spielbergs amistad based true story group africans board slave ship captured taken america legal disp  
993 capsule style heist film set present robert deniro stars wants retire form crime takes one last job request  
994 zero effect gets title main character daryl zero bill pullman although dont understand truly means last line  
995 asked see movie friend initial reaction hugh grant perhaps wrong harsh street hooker could picture romantic  
996 damn trailers advertising film reveals far much contents would glued sand film retains value thanks excellen  
997 youve got mail works alot better deserves order make film success cast two extremely popular attractive star  
998 trekkies roger energetic hilarious documentary brings viewers world star trek conventions beauty film good o  
999 set wild west carry around arrival rumpo kid sidney james cronies city dealings summary shootings judge burk  
1000 anxious see long time friend mine recommended crush neve campbell wanted prove shes hot thinks proved right

# Main Execution

- The predefined vocabulary is loaded from the file '*vocab.txt*'.
- All negative reviews (from '*txt\_sentoken/neg*' directory) are processed and saved to '*negative.txt*'.
- Similarly, all positive reviews (from '*txt\_sentoken/pos*' directory) are processed and saved to '*positive.txt*'.

```
# load doc into memory
def load_doc(filename):
    # open the file as read only
    file = open(filename, 'r')
    # read all text
    text = file.read()
    # close the file
    file.close()
    return text
```

```
# load vocabulary
vocab_filename = 'vocab.txt'
vocab = load_doc(vocab_filename)
vocab = vocab.split()
vocab = set(vocab)

# prepare negative reviews
negative_lines = process_docs('txt_sentoken/neg', vocab)
save_list(negative_lines, 'negative.txt')

# prepare positive reviews
positive_lines = process_docs('txt_sentoken/pos', vocab)
save_list(positive_lines, 'positive.txt')
```

# Processing All Documents in a Directory (*process\_docs*)

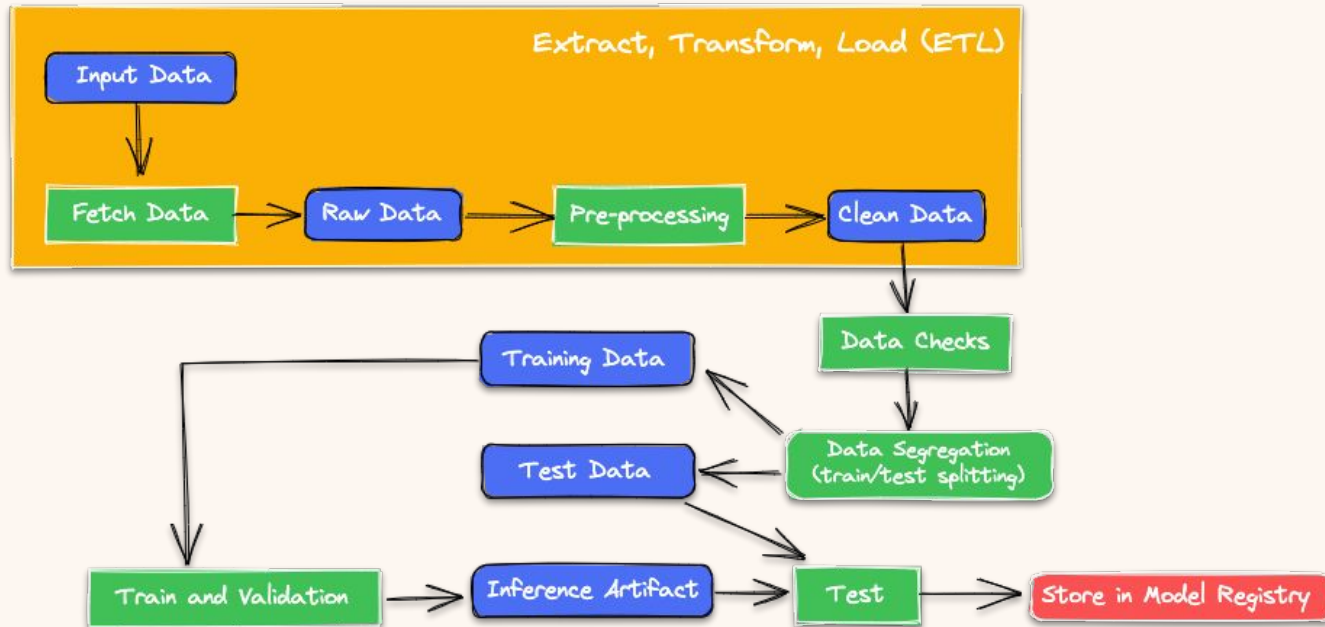
- This function processes all ".txt" files in a given directory.
- It loads each file, cleans it, and keeps words from the predefined vocabulary.
- The cleaned content of each document is added to a list.

```
# load all docs in a directory
def process_docs(directory, vocab):
    lines = list()
    # walk through all files in the folder
    for filename in listdir(directory):
        # skip files that do not have the right extension
        if not filename.endswith(".txt"):
            continue
        # create the full path of the file to open
        path = directory + '/' + filename
        # load and clean the doc
        line = doc_to_line(path, vocab)
        # add to list
        lines.append(line)
    return lines
```

```
# load doc, clean and return line of tokens
def doc_to_line(filename, vocab):
    # load the doc
    doc = load_doc(filename)
    # clean doc
    tokens = clean_doc(doc)
    # filter by vocab
    tokens = [w for w in tokens if w in vocab]
    return ' '.join(tokens)
```



When we look at the pipeline, we realize that what was done to clean the dataset has a serious flaw, what would it be?



Cont.