

Blueprint: Next-Gen Enterprise RAG & LLM

Nvidia PDFs Use Case

Vincent Granville, Ph.D.
vincentg@MLTechniques.com
www.MLTechniques.com
Version 1.0, January 2025

Contents

1	Introduction	1
2	From Corpus to xLLM Format	2
2.1	Hierarchical chunking with multi-index	2
2.2	Duplicate and overlapping chunks	3
2.3	PDF parser and data	3
3	Backend architecture	9
3.1	Backend tables	9
3.2	Backend parameters	10
3.3	Backend stemming	10
3.4	Adding tags, categories, or agents	11
4	Frontend architecture	12
4.1	Frontend tables	12
4.2	Frontend parameters	12
4.3	Frontend stemming	13
4.4	Scoring engine	14
4.5	Advanced user interface with-real time fine-tuning	14
4.6	Testing and Evaluation	15
5	Python code	16
5.1	Main program	16
5.2	Internal library	23
	References	27
	Index	29

1 Introduction

In this paper, I explain how to build an enterprise LLM and RAG system from scratch, without training, GPU, neural network or transformer, while delivering more accurate and exhaustive results at a fraction of the cost. It includes the full code with documentation, a deep dive on each component, the advanced UI with real-time fine-tuning, relevancy scores and structured output, the efficient in-memory database architecture, as well as front-end and back-end parameters. Each with detailed description and latest updates regularly added.

To illustrate the methodology, I use a Nvidia repository consisting of financial documents stored as PDFs. Also featured here: deep PDF retrieval including tables, images, bullet lists and font types, multi-index hierarchical and contextual chunking, as well as categories or agents assignment at the backend (chunking) level. Reliance on external Python libraries is minimum and complemented by enhancements as needed, to avoid drawbacks that can cause hallucinations. You will find here the secret sauce and innovative thinking behind LLM 2.0 and Enterprise xLLM in particular, based on many years of experience and technological vision.

The focus is on the core architecture. Specialized applications such as LLM for [cataloging](#), [text clustering](#) and predictive analytics, [sub-LLMs](#) and [LLM router](#), DNA sequences and tabular data [synthetization](#), API calls for problem solving and code generation, are covered in books [\[1, 2\]](#). In particular, Part 1 and Chapter 10 in [\[1\]](#) are most relevant to the material presented here. The most recent version of this article is available [here](#).

2 From Corpus to xLLM Format

One of the key features of this pre-processing step is to extract contextual elements when turning the corpus data into the input format required by xLLM, our proprietary model. In the case of the Wolfram corpus consisting of 15k webpages and 5k categories, the goal is to recover and use the embedded taxonomy and other contextual elements, rather than building a knowledge graph on top of the website, from scratch. For the Inbev corporate corpus, contextual elements present in the corpus include breadcrumbs, categories, tags, titles, and so on. For both examples, see details in my books [1, 2].

In the example discussed here (Nvidia PDF repository) the context includes document title, page title and number, text in big or bold font, sub-lists in bullet lists, and more. In addition, to each chunk, I assign extra contextual elements such as tag, category, or agent. For instance, the ‘table’ or ‘image’ agent indicates that the chunk contains tables or images. I describe **chunking** in section 2.1 and **agents** in section 3.4. An LLM dealing not only with text but also images and other media, is called **multimodal**.

2.1 Hierarchical chunking with multi-index

Chunks are text entities delimited by separators, such as sentences, paragraphs, or pages. In this version, I use two levels: paragraphs and pages. For details, see Chapter 10 in [1], also discussing table, font size, bullet lists and sublists detection. A **multi-index** is attached to each chunk, for easy retrieval and browsing. See the two sample chunks in this section. The index, referred to as ID in the Python code, is A6X2 for the first chunk and B3051X2 for the second one. The letter A indicates a parent chunk, while B stands for a **sub-chunk**. The strings X2, X1 indicate that the chunks are respectively in the 2nd and 1st PDF document. Images and tables are stored separately, with an index that includes the document and page number, respectively doc_ID and pn in the code.

```
A6X2~~{'title_text': '', 'description_text': '| "NVIDIA Builds Out Its|Omniverse  
Ecosystem|to Support the|Automotive Metaverse"|. |-- SiliconANGLE|. |NVIDIA  
Omniverse is an ideal tool for an|industrial world seeking to  
digitalize.|Omniverse can simulate the best|possible layouts before the first  
brick|is placed. The $3 trillion automotive|industry is modernizing all its  
processes|to take advantage of computing and|AI. BMW Group uses Omniverse  
to|build a whole factory digital twin before|constructing it physically.  
Mercedes-|Benz is using NVIDIA DRIVE IX on|Omniverse to design and simulate  
its|integrated cabins and electronics.|', 'doc_ID': 2, 'pn': 10}  
  
B3051X1~~{'title_text': '', 'description_text': '| $|1,990|$|847|$|10,896|', 'type':  
'Data', 'doc_ID': 1, 'pn': 170, 'block_ID': 96, 'item_ID': -1, 'sub_ID': -1,  
'fs': 8.7, 'fc': 0, 'ft': 'DIN_Next_LT_Pro_Light'}
```

Besides images, the xLLM input format is plain ASCII text with a structure similar to JSON, with special or accented characters being encoded. Each row corresponds to a chunk or sub-chunk, and is stored as a hash. The hash key indicates the type of information, for instance:

- page title: title_text,
- full text: description_text,
- page number pn, document ID: doc_ID,
- type of content: type (the value Data indicates data in chunk B3051X1),
- list of tags, agents, categories, parent categories, images, tables
- font size: fs, font color: fc

and so on. Content types include data, bullet list, sub-list, note, and more. Another useful hash key is URLs, in particular, the URL pointing to the specific PDF page where the chunk is located. For instance, if doc.pdf is on example.com, then example.com/doc.pdf#page=5 points to page 5 in the PDF. Note that the the number of hash keys in a chunk is not fixed. The full xLLM input file (the chunks stored as key-value pairs) is available [here](#).

The backend table Index_to_IDs generated by the main program in section 5.1, connects each parent chunk to its sub-chunks. It has one row per parent chunk. In the example below featuring one row, the parent chunk is indexed as A18X1, while such-chunks have an index starting with letter B. The vector at the beginning indicates the document ID and page number. The numbers in black represent the size of each chunk. You can download the full table is available [here](#). Usually, the parent chunk is larger than its sub-chunks. The exception

is when the parent chunk is very small: sub-chunks have more overhead that can make them bigger than the parent chunk, but this is usually an indication of poor chunking.

```
(1, 22) {'B192X1': 205, 'B193X1': 1030, 'B194X1': 255, 'B195X1': 230, 'B197X1': 184,
        'B198X1': 201, 'B199X1': 185, 'B200X1': 201, 'B201X1': 185, 'B202X1': 202,
        'B203X1': 186, 'B204X1': 244, 'B205X1': 365, 'B206X1': 429, 'B207X1': 468,
        'B208X1': 575, 'B209X1': 174, 'A18X1': 2396}
```

2.2 Duplicate and overlapping chunks

Chunks in different PDFs may be identical. For instance, the tables of contents are identical in the 2022 and 2023 financial reports (two separate PDFs), as it based on the same template. Whether or not keeping them both is a business decision. Also, within a same PDF, some parent chunks may contain only one sub-chunk. In this case, both are identical due to poor chunking.

In the current version, sub-chunks do not overlap, and parent chunks do not overlap either. However, it is a good practice to allow for **overlapping chunks**. This is easy to achieve since chunks are sequentially generated. For instance, it makes sense to group sub-chunks into bigger chunks, with the last sub-chunk of a big chunk being the first sub-chunk of the next big chunk. The reason for doing this is because a token in a sub-chunk may be related to another token in the next sub-chunk. Parent chunks partially address this issue, but there is room for improvement.

Another question is why storing both sub-chunks and parent chunks, when sub-chunks alone are enough as they can be used to reconstruct the parents chunks. The reason is due to the relatively small size of enterprise corpuses. If there is enough space and memory to store both, it is more efficient to keep them both despite doubling the memory requirements. Note that the space needed in backend tables to store all the elements of a chunk grows faster than linearly with the size of the chunk: in particular, the number of token pairs consisting of tokens found in a same chunk (used to build embeddings) grows quadratically with the size of the chunk. To avoid this problem, tokens that are not close enough to each other cannot form a token pair; the maximum distance allowed is determined by the backend parameter `maxDist`, see section 3.2.

2.3 PDF parser and data

The Python code is similar to the older version documented in Chapter 10 in [1]. The main differences are:

- Dealing with multiple PDFs. See lines 292–302.
- Detection of duplicate chunks (when a parent chunk only has one sub-chunk), see lines 317–318.
- Multi-index with specific codes to distinguish chunks from sub-chunks, see lines 345 and 380.

The code is on Google Drive, [here](#). The full xLLM output file is available [here](#). The three PDF documents used as input in the code can be found respectively [here](#), [here](#), and [here](#). In the Python code, lines 85–91 deal with special characters. Pages of each PDF are saved as images with the proper index in lines 396–422. Note the image detection mechanism in lines 265–274, table detection using the **PymuPDF** library in line 169–179, as well as home-made bullet list, sub-list and table detection for tables missed by PymuPDF. For detailed description, see Chapter 10 in [1]. **PDF_Chunking_Nvidia.py**

```
1 # Input PDF for this script: https://drive.google.com/file/d/1Daa9oZJm4-b6NqUsVGxK2euemcFnf8jH/
2
3 # pymupdf.readthedocs.io/en/latest/page.html#Page.find_tables
4 # stackoverflow.com/questions/56155676/how-do-i-extract-a-table-from-a-pdf-file-using-pymupdf
5
6 import fitz # PyMuPDF
7
8
9 def update_item_ID(k, entity_idx, type, table_ID):
10
11     idx = entity_idx[k]
12     if type == 'Data':
13         # table: data row
14         flag = 'TD' #
15     elif type == 'Note':
16         # table: labels
17         flag = 'TL'
18     idx_list = list(idx)
19     idx_list[1] = flag + str(table_ID)
20     entity_idx[k] = tuple(idx_list)
21
```

```

22     return(entity_idx)
23
24
25 def detect_table(xLLM_entity):
26
27     # detect and flag simple pseudo-tables
28
29     entity_txt = xLLM_entity[0]
30     entity_type = xLLM_entity[1]
31     entity_idx = xLLM_entity[2]
32     table_ID = -1
33     table_flag = False
34
35     for k in range(1, len(entity_type)):
36
37         type = entity_type[k]
38         text = entity_txt[k]
39         old_text = entity_txt[k-1]
40         old_type = entity_type[k-1]
41
42         if ( (
43             (type == 'Data' and old_type == 'Note') or
44             (type == 'Note' and old_type == 'Data') or
45             (type == 'Data' and old_type == 'Data')
46         )
47             and old_text.count('|') == text.count('|')
48             and text.count('|') > 2
49         ):
50             print("detected table", table_ID + 1)
51             if not table_flag:
52                 table_ID += 1
53                 table_flag = True
54             idx = entity_idx[k]
55             old_idx = entity_idx[k-1]
56
57             # update item_ID (idx[1] and old_idx[1]) in current and previous row
58             # item_ID starts with letter D if data, or N if labels
59
60             update_item_ID(k, entity_idx, type, table_ID)
61             update_item_ID(k-1, entity_idx, old_type, table_ID)
62
63         else:
64             #table_ID = -1
65             table_flag = False
66
67     return(xLLM_entity)
68
69
70 def cprint_page(xLLM_entity, hash_chunks, OUT):
71
72     entity_txt = xLLM_entity[0]
73     entity_type = xLLM_entity[1]
74     entity_idx = xLLM_entity[2]
75
76     for k in range(len(entity_type)):
77
78         type = entity_type[k]
79         text = entity_txt[k]
80         text = text.strip()
81         text = text.replace(" ", " ")
82         text = text.replace(" |", "|")
83         text = text.replace("| ", "|")
84         text = text.replace("||", "|")
85         text = text.encode('unicode-escape').decode('ascii')
86         text = text.replace('\u2022', '.')
87         text = text.replace('\u2013', '--')
88         text = text.replace('\u2014', '--')
89         text = text.replace('\u2019', '"')
90         text = text.replace('\u201c', '"')
91         text = text.replace('\u201d', '"')
92
93         idx = entity_idx[k]
94         doc_ID = idx[3]
95         block_ID = idx[0]
96         item_ID = idx[1]
97         sub_ID = idx[2]

```

```

98     pn = idx[4] # page number
99     fs = idx[5] # font size
100    fc = idx[6] # font color
101    ft = idx[7] # font typeface
102    #- print(k, type, idx, text)
103    OUT.write(f"{type:<8}{block_ID:>3}{item_ID:>5}{sub_ID:>3}{pn:>3}"
104              f"{fs:>5}{fc:>9} {ft:<20}{text:<80}\n")
105    hash_chunks[idx] = (type, text) # ignore if type == Data ??
106
107    OUT.write("\n")
108    return(hash_chunks)
109
110
111 def update_page(text, type, entity, idx):
112
113     entity_txt = entity[0]
114     entity_type = entity[1]
115     entity_idx = entity[2]
116     block_ID = idx[0]
117     item_ID = idx[1]
118     sub_ID = idx[2]
119     k = len(entity_txt)
120     if k > 0:
121         old_type = entity_type[k-1]
122         old_idx = entity_idx[k-1]
123         old_block_ID = old_idx[0]
124         old_item_ID = old_idx[1]
125         old_sub_ID = old_idx[2]
126     else:
127         old_type = ""
128         old_block_ID = ""
129         old_item_ID = ""
130         old_sub_ID = ""
131     if type in ('Note', 'Data'):
132         sep = "|"
133     else:
134         sep = " "
135
136     if (type == old_type and block_ID == old_block_ID
137         and item_ID == old_item_ID and sub_ID == old_sub_ID):
138         new_text = entity_txt[k-1] + text + sep
139         entity_txt[k-1] = new_text
140     else:
141         entity_txt.append(sep + text + sep)
142         entity_type.append(type)
143         entity_idx.append(idx)
144     return(entity)
145
146
147 def convert_pdf_to_json(doc_ID, filename):
148
149     pdf_path = filename + '.pdf' # input PDF
150     json_path = filename + '.json' # PDF turned to Json (unused)
151     text_path = filename + '.txt' # PDF turned to text (unused, for debugging)
152
153     OUT = open(text_path, "wt", encoding="utf-8")
154     hash_chunks = {}
155
156     # Open the PDF file
157     pdf_document = fitz.open(pdf_path)
158     content = ""
159
160     # Iterate through the pages
161     for page_num in range(len(pdf_document)):
162         OUT.write("\n-----\n")
163         OUT.write("Processing page " + str(page_num) + "\n\n")
164         print("Page:", page_num)
165         page = pdf_document.load_page(page_num)
166
167         text_data = page.get_text("dict") # also extract as "json" to get tokens in green font
168
169         tabs = page.find_tables()
170         for tabs_index, tab in enumerate(tabs):
171             # iterate over all tables
172             index = (page_num, tabs_index)
173             table_data = tab.extract() # extracting tabs[i], the i-th table in this page

```

```

174         if len(table_data) > 0:
175             # if not, ignore this table (note the important parameter threshold here)
176             OUT.write("Table " + str(index) + ":\n")
177             for row in table_data:
178                 OUT.write(str(row) + "\n")
179             OUT.write("\n")
180
181     itemize = False
182     item_ID = -1
183     sub_ID = -1
184     block_ID = -1
185     item = ""
186     fsm = -1 # top level font size in bullet list
187     fst = 64 # min title font size (top parameter)
188     title = ""
189     notes = ""
190     old_block_number = -1
191     old_font_size = -1
192     entity_txt = []
193     entity_idx = []
194     entity_type = []
195     type = ""
196     old_text = ""
197     old_type = ""
198
199     for block in text_data["blocks"]:
200         if block["type"] == 0: # Text block
201             block_number = block["number"]
202             for line in block["lines"]:
203                 for span in line["spans"]:
204
205                     text = span["text"]
206                     font_name = span["font"]
207                     font_size = span["size"]
208                     font_size = round(font_size,1)
209                     font_color = span["color"]
210
211                     if font_size > fst:
212                         type = 'Title'
213                     elif ord(text[0]) == 8226:
214                         itemize = True
215                         if fsm == -1:
216                             fsm = font_size
217                         if font_size > 0.98 * fsm:
218                             item_ID += 1
219                             type = 'List'
220                         else:
221                             sub_ID += 1
222                             type = 'SubList'
223                     elif itemize:
224                         #- itemize = ((0.99 < font_size/old_font_size < 1.01) and
225                         #-             (not text[0].isupper() or ord(old_text[0]) == 8226))
226                         itemize = ((0.99 < font_size/old_font_size < 1.01) or
227                                   (ord(old_text[0]) == 8226))
228
229                     if not itemize:
230                         item_ID = -1
231                         sub_ID = -1
232                         block_ID += 1
233                         type = 'Note'
234                     else:
235                         if not text[0].isdigit() and text[0] not in ('$','+', '-'):
236                             type = 'Note'
237                         #- elif block_number != old_block_number:
238                         else:
239                             type = 'Data'
240
241                     if block_ID == -1:
242                         block_ID += 1
243                     elif ((type not in (old_type, 'List', 'SubList')) or
244                           (not (0.99 < font_size/old_font_size < 1.01))):
245                         if (old_text != "" and ord(old_text[0]) != 8226 and
246                             type not in ('List', 'SubList')):
247                             block_ID += 1
248
249     idx = (block_ID, item_ID, sub_ID, doc_ID, page_num, font_size,

```

```

250         font_color, font_name, block_number)
251     entity = (entity_txt, entity_type, entity_idx)
252     update_page(text, type, entity, idx)
253
254     old_font_size = font_size
255     old_text = text
256     old_type = type
257
258     old_block_number = block_number
259
260     entity = detect_table(entity)
261     cprint_page(entity, hash_chunks, OUT)
262
263
264     image_list = page.get_images()
265     for image_index, img in enumerate(image_list, start=1):
266         xref = img[0]
267         base_image = pdf_document.extract_image(xref)
268         image_bytes = base_image["image"]
269         size = len(image_bytes)
270         ext = base_image["ext"]
271         index = (page_num, image_index)
272         OUT.write(f"Image {str(index):<8}{str(ext):>5} size = {str(size):>6}\n")
273         #- with open(f"image_{page_num + 1}_{image_index}.{ext}", "wb") as image_file:
274         #-     image_file.write(image_bytes)
275
276     content += "__P" + str(page_num) + "\n" + page.get_text("json") # "text", "html", "json"
277
278     OUT.close()
279     with open(json_path, 'w', encoding='utf-8') as json_file:
280         json_file.write(content)
281     return(hash_chunks)
282
283
284 # --- Main ---
285
286 def get_value(hash, key):
287     if key in hash:
288         return(hash[key])
289     else:
290         return('')
291
292 filenames = (
293     'nvda-f2q24-investor-presentation-final-1',
294     '2022-Annual-Review',
295     '2023-Annual-Report-1',
296 )
297
298 outfilename = "repository_Nvidia_X.txt" # xLLM input corpus
299 OUT_xllm = open(outfilename, "wt")
300
301
302 #-- loop over PDFs
303
304 entity_list = ()
305
306 for doc_ID in range(len(filenames)):
307
308     filename = filenames[doc_ID]
309     hash_chunks = convert_pdf_to_json(doc_ID, filename)
310
311     #-- Add small chunks to xLLM input
312
313     hash_chunkTitle = {}
314     for idx in hash_chunks:
315         content = hash_chunks[idx]
316         # next 2 code lines to track duplicate chunks
317         if content[1] not in entity_list:
318             entity_list = (*entity_list, content[1])
319         if content[0] == "Title":
320             hash_chunkTitle[(idx[3], idx[4])] = content[1]
321
322     uID = 0
323     hash_entities = {}
324
325     for idx in hash_chunks:

```

```

326     content = hash_chunks[idx]
327     title = get_value(hash_chunkTitle, (idx[3],idx[4]))
328     hash_entities[uID] = {
329         'title_text': title,
330         'description_text': content[1],
331         'type': content[0],
332         'doc_ID': idx[3],
333         'pn': idx[4], # page number
334         'block_ID': idx[0],
335         'item_ID': idx[1],
336         'sub_ID': idx[2],
337         'fs': idx[5], # font size
338         'fc': idx[6], # font color
339         'ft': idx[7], # font typeface
340     }
341     uID += 1
342
343 for ID in hash_entities:
344     # print(ID, hash_entities[ID])
345     findex = 'B' + str(ID) + 'X' + str(doc_ID)
346     OUT_xllm.write(findex + "~~" + str(hash_entities[ID]) + "\n")
347
348 #-- Add big chunks to xLLM input
349
350 hash_bigEntities = {}
351
352 for idx in hash_chunks:
353     content = hash_chunks[idx]
354     title = get_value(hash_chunkTitle, (idx[3],idx[4]))
355     pn = idx[4], # page number
356     doc_ID = idx[3]
357     ID = (pn, doc_ID)
358     text = content[1]
359     type = content[0]
360     if ID in hash_bigEntities:
361         local_hash = hash_bigEntities[ID]
362         local_hash['description_text'] += ". " + text
363         if local_hash['title_text'] == '' and title != '':
364             local_hash['title_text'] = title
365         hash_bigEntities[ID] = local_hash
366     else:
367         hash_bigEntities[ID] = {
368             'title_text': title,
369             'description_text': content[1],
370             'doc_ID': idx[3],
371             'pn': idx[4], # page number
372         }
373
374 vID = 0
375 for ID in hash_bigEntities:
376     # print(ID, hash_bigEntities[ID])
377     local_hash = hash_bigEntities[ID]
378     text = local_hash['description_text']
379     if text not in entity_list:
380         findex = 'A' + str(vID) + 'X' + str(doc_ID)
381         OUT_xllm.write(findex + "~~" + str(hash_bigEntities[ID]) + "\n")
382     else:
383         # this big chunk is identical to a small chunk
384         # maybe because small chunks are not granular enough
385         # maybe due to identical big chunks found in multiple PDFs
386         # ignore to avoid duplicate chunks, to save space
387         doc_ID = local_hash['doc_ID']
388         pn = local_hash['pn']
389         print("Ignored: Doc_ID=%3d pn=%3d vID=%4d" % (doc_ID, pn, vID))
390         pass
391     vID += 1
392
393 OUT_xllm.close()
394
395 # --- PDF to Images ---
396
397 def PDF_to_PNG(doc_ID):
398
399     from PIL import Image
400
401

```



```

402 pdf_path = filenames[doc_ID] + '.pdf'
403 pdf_document = fitz.open(pdf_path)
404 zoom = 2 # to increase the resolution
405 mat = fitz.Matrix(zoom, zoom)
406
407 for page_num in range(len(pdf_document)):
408     page = pdf_document.load_page(page_num)
409     pix = page.get_pixmap(matrix = mat) # or (dpi = 300)
410     img = Image.frombytes("RGB", [pix.width, pix.height], pix.samples)
411     filename = "PDF" + str(doc_ID) + "_" + str(page_num) + '.png'
412     img.save(filename)
413     print('Converting PDFs to Image ... ' + filename)
414
415     return()
416
417 # set save_PNG = True to save PDF slides as PNG
418
419 save_PNG = False
420 if save_PNG:
421     for doc_ID in (0, ):
422         PDF_to_PNG(doc_ID)

```

3 Backend architecture

The architecture is essentially the same as described in Part 1 in [1], relying on in-memory nested hashes for **backend** tables, to increase speed by orders of magnitude. These tables are generated on the first run or when you update the backend parameters, and then saved as text files. It takes a few seconds to build them. On subsequent runs, they are imported for even faster processing. By backend, I mean corpus processing, while **frontend** deals with prompt processing.

There are several components that differentiate **xLLM** – our architecture – from standard LLMs, such as no transformer, no **prompt engineering**, no training, no weight, yet increased speed, accuracy, and exhaustivity. Also, **explainable AI** (intuitive parameters), much lower cost, **reproducibility**, and almost no reliance on external libraries or APIs, making the system more secure. The version described here is **Enterprise xLLM**.

Other improvements include several token types, long multi-tokens, **variable length embeddings**, smart relevancy scores and ranking, generic **PMI** (pointwise mutual information) instead of cosine similarity, and hierarchical chunking with multi-index. The concepts new to this version are described in section 2, 3.3, 3.4, and 4. Here I describe the main token types. Tokens consist of one or multiple words. In the latter case, words are separated by '~'. A word such as 'San Francisco' is a single token, though 'San' and 'Francisco' may also be single tokens, offering some redundancy. The following token types encompass both single and **multi-tokens**.

- **s-tokens** come from the stemming process: if 'gaming' is in the prompt, the words 'game' and 'games' will also be searched for. Because they are different from the original word, they are flagged as s-tokens, and given a lower weight in the relevancy algorithm. Soon to be implemented.
- **g-tokens** are tokens found in the contextual fields attached to chunks, such as agent, tag, title, or category. Named gword in the code, the prefix 'g' indicates that they are connected to the **knowledge graph**. They are given a higher weight in the relevancy algorithm.
- **c-tokens** are multi-tokens consisting of words that are not adjacent in the corpus. I do not use them here.

3.1 Backend tables

There are four types of backend tables:

- Auxiliary tables for stemming or un-stemming (hash_stem used in the backend, hash_unstem used in the frontend but built in the backend) and the stopwords list. Also in the same category, Index_to_ID maps chunk IDs to the corresponding document and page numbers; these two numbers constitute the table key.
- Indexed by a **chunk**. They start with 'ID' in the table name, where ID is the **multi-index** attached to the chunk in question. See also section 2.1.
- Indexed by a **multi-token**: dictionary, hash_agents, hash_ID, embeddings, sorted_ngrams, as well as those where the table name starts with hash_context. The later deals with the contextual fields such as title, tag, or category, as opposed to the actual content attached to a chunk.
- Indexed by a pair of related multi-tokens: sorted_ngrams. This table is used to create embeddings, to suggest related prompts to the user. It requires the large temporary hash_pairs table; the latter is not

saved or loaded when reading tables from text files, as it is used in the building process only, and never accessed from the frontend. The `ctokens` table, a contextual version of `hash_pairs`, is not used here.

Backend tables are loaded in memory. Many are stored as **nested hashes**. They are saved as text files. To read them, see the `read_backend_tables` function in the code in section 5.1, itself relying on sub-functions from our internal library. That library is imported as `exllm`, see line 4 in the same code. Backend tables are defined in lines 12–38. Also, you can find them on GitHub, [here](#). Each row correspond to a specific key in the table.

Reading these tables is greatly facilitated using the `eval` Python function that turns strings into nested hashes and can detect integers and floats. Thanks to it, I was able to remove and simplify a number of reading functions in our internal library; the price to pay is a slight reduction in speed. To build the tables (as opposed to reading them), set `build_tables` to `True`. When running the code, you will see the tables – whether built or read – displayed on the screen as in Figure 1, with the number of rows per table, and the table type.

```

Saving table      dictionary 170146 elts <class 'dict'>
Saving table      hash_context1 170146 elts <class 'dict'>
Saving table      hash_context2 170146 elts <class 'dict'>
Saving table      hash_context3 170146 elts <class 'dict'>
Saving table      hash_context4 170146 elts <class 'dict'>
Saving table      hash_context5 170146 elts <class 'dict'>
Saving table      hash_ID 170163 elts <class 'dict'>
Saving table      hash_agents 170146 elts <class 'dict'>
Saving table      ID_to_content 8709 elts <class 'dict'>
Saving table      ID_to_agents 40 elts <class 'dict'>
Saving table      ID_size 9627 elts <class 'dict'>
Saving table      ID_to_index 8709 elts <class 'dict'>
Saving table      Index_to_IDs 430 elts <class 'dict'>
Saving table      stopwords 75 elts <class 'tuple'>
Saving table      hash_stem 2295 elts <class 'dict'>
Saving table      hash_unstem 6052 elts <class 'dict'>
Saving table      embeddings 38001 elts <class 'dict'>
Saving table      sorted_ngrams 166291 elts <class 'dict'>

```

Figure 1: Backend tables, with type and number of rows

3.2 Backend parameters

Backend parameter values are set in the `build_backend_tables` function. See code in section 5.1, also shown below for convenience.

```

backendParams = {
    'max_multitoken': 4, # max. consecutive terms per multi-token
    'maxDist' : 3,      # max. distance between 2 multi-tokens to link them in hash_pairs
    'create_hpairs': True, # required for embeddings, takes time and memory
    'create_ctokens': False, # needs create_hpairs set to TRUE for activation
    'use_stem': False,      # backend stemming
    'extraWeights':         # default weight is 1
        {
            'description': 0.0,
            'category': 0.0, # try 0.5
            'tag_list': 0.0, # try 0.5
            'title': 0.0, # try 0.5
            'meta': 0.0 # try 0.5
        }
}

```

If you set `create_hpairs` to `False`, embeddings won't be created. It saves some time and space, but prevents the user from seeing related keywords that could be used for extra prompting. The `extraWeights` values allow you to boost the importance of multi-tokens found in **contextual elements**, compared to those in the standard text. Finally, a higher `maxDist` increases the width of the **contextual window** for embeddings, but also leads to more memory usage, and bigger `hash_pairs` and `embeddings` tables especially for large chunks. It leads to more results potentially shown to the user, though **distillation** and the **scoring engine**, controlled by the frontend parameters, will only pick up only the most relevant ones.

3.3 Backend stemming

With `'use_stem': True` in `backendParams`, words in the corpus are replaced by their most common alternative that is also present in the corpus. Table 1 shows the original word in column `token1`, with its substitution

in column `token2`. Columns n_1, n_2 are the number of occurrences in the corpus, respectively for the original word, and its substitution. It allows you to easily retrieve and group the different versions of a corpus word. It also results in smaller backend tables and slightly faster processing. But it sometimes produce inaccurate answers when the **stemming** is too aggressive. The default is `'use_stem': False`.

n_1	n_2	token ₁	token ₂
8	42	region	regions
12	42	regional	regions
8	188	register	registered
26	70	registrants	registrant
12	70	registration	registrant
85	182	regulation	regulations
2	182	regulator	regulations
16	182	regulators	regulations
8	14	reimburse	reimbursement
4	14	reimbursed	reimbursement
4	14	reimbursements	reimbursement
2	12	reincorporation	reincorporated
2	6	reinforce	reinforced
4	6	reinvent	reinvented
2	6	reinventing	reinvented
2	6	reinvents	reinvented
4	12	reinvest	reinvested
8	12	reinvestment	reinvested

Table 1: Backend table `hash_stem` (extract)

The code below, in the `update_dict` function in the `xllm-enterprise-nvidia-util` library, does the substitution. The algorithm behind the scenes is explained in section 4.3. Note that n_1, n_2 in Table 1 are usually even integers due to double counting: each token appears both in a **sub-chunk** and in its **chunking**, unless the parent has only one child.

```
for word in words:
    if use_stem and word in hash_stem:
        word = hash_stem[word]
```

3.4 Adding tags, categories, or agents

There are two flavors of **agents**. In both cases, the goal is to assess the user intent and serve results accordingly. **Retrieval agents** focus on detecting the type of content the user is looking for: best practices, definitions, templates, datasets, code, spreadsheets, images, tables, and so on. **Action agents** perform executive functions, such as solving a math problem, generating code, generating synthetic data, auto-tagging, auto-indexing, cataloging, summarizing, predictive analytics, or text clustering.

Action agents rely on API calling external apps such as Wolfram to solve math problems. However, xLLM has high quality internal apps for a number of tasks, such as cataloging, text clustering, predictive analytics on text data, and tabular data synthetization. In contrast to standard LLMs that try to assign agents to prompts, xLLM offers the user the option to select retrieval agents listed on the command menu. Also, agents are assigned to chunks along with other contextual elements such as titles, tags, or categories. Unlike titles, agents are assigned post-crawling based on business needs and corpus content. In the Nvidia use case, images and tables are two possible values for the agent. To create a meaningful list of agents, you need to look at top multi-tokens found in the dictionary backend table. Tags and categories are built using the same principle.

In the current version of xLLM, agents are detected and assigned to chunks in the `build_backend_tables` function in `xllm-enterprise-nvidia-dev.py`. Look for the keyword 'agent' in the code in question. The plan is to have them assigned earlier when building the xLLM file with the PDF parser. Either way, they end up in the `ID_to_agents` backend table. In this table, the key is a chunk represented by its ID (multi-index), and the value is the list of associated agents. Chunks with no agent are not included in that table. This allows for easy agent retrieval in the frontend. See `agent_map` for the current list of agents.

ID	ID_agents
B2X0	('income', 'tax', 'cash', 'products', 'data', 'accelerated computing', 'data center', 'non-gaap')
A154X1	('stock', 'income', 'tax', 'cash', 'revenue', 'common stock', 'financial statements')
A175X1	('stock', 'cash', 'data', 'assets', 'financial statements', 'restricted stock', 'data center')
A179X1	('cash', 'securities', 'assets', 'financial statements')
A182X1	('income', 'cash', 'equity', 'securities', 'assets', 'financial statements')
A186X1	('products', 'assets', 'financial statements')
B2250X2	('stock', 'financial statements')
B2749X2	('stock', 'income', 'tax', 'cash', 'directors', 'common stock', 'securities', 'financial statements')

Figure 2: Sample prompt results, agents view corresponding to summary view in Figure 3

Figure 2 shows the **agents view**, that is, the agent list attached to each chunk relevant to the prompt. In this example, the prompt is ‘financial statements 2024’, and chunks are represented by their ID (a **multi-index**). I also discuss this prompt in section 4.5, featuring the summary and embeddings views. These three views are part of the prompt results. Here, we are dealing with **retrieval agents**. All of them are **soft agents**, defined as keywords found in the chunk content, as opposed to **hard agents**, defined as agents attached to the contextual fields of a chunk such as tags, categories, or agents. Hard agents have names that start with a capital letter, by contrast to soft agents. Note that soft retrieval agents play the same role as pre-assigned tags. They are useful to decide whether or not you want to see the full content of a specific chunk, before clicking on it to get the detailed view.

4 Frontend architecture

The **frontend** deals with how to connect prompt tokens to backend tables, what results to display to the user, in what order, and how: structured output with concise but exhaustive, original results not reworded differently, by contrast to standard LLMs. The main components are **unstemming**, sorted ***n*-grams**, frontend **distillation**, and the **relevancy scores**. The UI allows for **real-time fine-tuning**. It also offers **agents** or categories to choose from, as well as entering multi-tokens or **negative keywords** in the prompt box. Finally, it displays **embeddings** to suggest related prompts, and the ability to retrieve full **chunks** from the corpus, including tables, URLs, and images, in just one click. It also displays relevancy scores attached to each chunk featured in the prompt results.

4.1 Frontend tables

The two main tables are `q_dictionary` and `q_embeddings`. They are local versions of the dictionary and embeddings backend tables. The prefix ‘q’ indicates that they are restricted to the data detected in the user query (the prompt). The `sorted_ngrams` (built in the backend) allows you to very efficiently check all 2^n combinations of frontend keywords A_1, \dots, A_n in the prompt, to identify the few ones found in the corpus, and stored as multi-tokens in the dictionary. If frontend stemming is activated, the backend table `hash_unstem` is also used.

Also, `ID_hash` is a local (prompt-related), transposed version of the backend table `hash_ID`. It is a multi-level key-value table, also called **nested hash**. The key is a corpus chunk represented by its multi-index. The value is a hash table (sub-hash) containing the multi-tokens found simultaneously in the corresponding chunk and in the prompt. The sub-key is a multi-token, and the sub-value is the multi-token multiplicity: the number of occurrences of the multi-token in question, in the parent chunk.

Finally, `ID_score` is built using `ID_hash`, and assigns a multivariate score measuring the relevancy between a specific chunk, and the prompt. The **score vector** is made of **sub-scores**, each measuring a particular aspect of relevancy: for instance, the presence of contextual multi-tokens matching those in the prompt, the number of tokens found simultaneously in a chunk and in the prompt, and so on. Each sub-score is sorted separately to produce sub-ranks, see `ID_score_ranked`. This simple hash table is indexed by a chunk ID; the value is a score vector. Then the final **rank** attached to a chunk, relative to a prompt, is a weighted sum of the sub-ranks.

4.2 Frontend parameters

The frontend parameter values are set in the `default_frontendParams` function in our internal library, see code in section 5.2. The list is much shorter compared to the previous version discussed in Part 1 in [1]. However, additional parameters will be added to better fine-tune the relevancy scores and distillation function. The PMI function is now moved to the backend, thus eliminating frontend PMI parameters.

To use frontend stemming, set `use_stem` to `True` and make sure any glitch from the Nltk Python library are addressed with the workaround discussed in section 4.3, to avoid low relevancy on occasions (or **hallucinations** in the case of standard LLMs). Note that `maxTokenCount` is attached to the frontend distillation function, activated when `distill` is set to `True`. The smaller the value, the stronger the impact. The parameter `beta`

is attached to the relevancy score. Several relevancy parameters will soon be added: what they will control is discussed in see section 4.4. The current frontend parameter list in the code below.

```
frontendParams = {
    'distill': False,
    'maxTokenCount': 1000, # ignore generic tokens if large enough
    'beta': 1.0, # used in text entity relevancy score, try 0.5
    'nresults': 20, # max number of chunks to show in results
    'use_stem': False,
}
```

In the end, the goal is to offer a small number of intuitive parameters based on [explainable AI](#), to make real-time fine-tuning easy for non-experts. As previously, there will be a catch-all parameter set for the expert user who wants to see the full output instead of the selection displayed on the screen. This helps understand the reason why some results may be missing depending on the parameter values. It is a valuable feature for [LLM debugging](#). Currently, for debugging, set `distill` and `use_stem` to `False`, and set `nresults` (the maximum number of chunks returned) to a very large value, say 1000.

4.3 Frontend stemming

The goal is to find in the corpus keywords related to those in the prompt. A simple example is ‘projection’ not found in the corpus, even though ‘projections’ is there. Also, if ‘gaming’ is a prompt token also present in the corpus, it makes sense to retrieve ‘games’ and ‘game’, not found in the prompt but present in the corpus. Yet a more complex example is ‘projection’ related to ‘prediction’ and ‘forecast’.

For [stemming](#), I use the PorterStemmer available in the [NLTK](#) Python library. See the `stem` function in our `xllm_enterprise_nvidia_util` library. It produces the `hash_unstem` backend table featured in Table 2, addressing a complex problem known as [unstemming](#). Despite being one of the least aggressive stemmers, it is still too aggressive for a system like ours that does not rely on [transformers](#). Some entries in `hash_unstem` must be broken down into smaller pieces to significantly improve performance. As an example, ‘project’ and ‘projection’ cannot have the same stem. The full `hash_unstem` table is in a zip archive on GitHub, [here](#).

Frontend stemming works as follows. If the word ‘games’ is in the prompt, xLLM looks for all alternate versions found in the corpus, such as ‘gaming’, yet giving a higher weight to the original word in the prompt. First, ‘games’ needs to be stemmed to ‘game’ using the PorterStemmer, then unstemmed to all variants using `hash_unstem[‘game’]`. This feature is activated if the frontend parameter `use_stem` is set to `True`.

stem	related corpus tokens
enterpris	enterprises, enterprise
race	racing, race, races
adopt	adopt, adoption, adopted, adopts, adopting
belief	belief, beliefs
global	global, globally
end	end, ends, ended, ending
game	gaming, games, game, gamings
return	returned, returns, return, returning
grow	growing, grows, grow
number	number, numbers
gpu	gpu, gpus
cloud	cloud, clouds
provid	providers, providing, provider, providence
deploy	deploying, deploy, deployed, deployment

Table 2: Backend table `hash_unstem` (extract)

If in addition, the backend parameter with same name `use_stem` is also set to `True`, then backend tables are smaller and retrieval is faster. It may result in lower accuracy and works with minimum stemming only. Note that stemming can be done semi-manually, by first sorting all single tokens (consisting of one word) in the dictionary backend table, and then group close neighbors while browsing the dictionary, avoiding false positives (grouping unrelated tokens) and false negatives (failure to identify two tokens as related).

4.4 Scoring engine

See high-level discussion in section 4.1, explaining the score vector, sub-scores, sorted scores and sub-ranks, and how the final rank attached to a corpus chunk is computed, relative to the prompt. This component is being enhanced, to better take into account **s.tokens** (synonyms to tokens in the prompt), **g.tokens** (found in contextual fields attached to chunks), single versus multi-tokens, token multiplicity, chunk size, token rarity, whether a token is first in the prompt or not, favoring recent content, and so on.

4.5 Advanced user interface with-real time fine-tuning

While most LLM and RAG systems offer nothing more than a basic search box to enter prompts, our platform comes with a UI offering several options. Prompt results are concise, well structured and allow the user to get a general overview before digging deeper via subsequent actions or prompts. See description in Chaper 3 in [1]. The main features include

- **Fine-tuning** frontend parameters in real time,
- The user can choose specific tags, categories (sub-LLMs) or agents on the menu when playing with the web API, or on the command line in the offline version,
- New to this version are the options to enter **negative keywords** and **multi-tokens** in the prompt query. For instance, a search that includes the multi-token ‘public~conference~call’ will only return chunks that contain all words glued together. If adding the negative keyword ‘!publication’ (with the exclamation point to indicate that the keyword is negative), it will ignore chunks that also include the word ‘publication’.
- Processing prompts in bulk with results saved in text file. Ability to re-use a previous prompt and modify them, or re-use a previous set of parameters, or the catch-all parameter to return everything found in the corpus and related to the prompt (useful for debugging purposes).
- Relevancy **scores** attached to each chunk displayed in the prompt results. The most relevant chunks are displayed each with a short summary, including associated contextual elements such as tags, categories, title, agents, chunk size, and relevant tokens matching those in the prompt. The user can click on specific chunks to see the full details (tables, images, full text and so on).

Figure 3 shows the **summary view** for the prompt ‘financial statements 2024’. It displays eight retrieved chunks: 5 sub-chunks with ID starting with A, and 3 parent chunks with ID starting with B. The columns ‘PDF’ and ‘pn’ respectively show the PDF identifier and page number within that document. Column ‘wRank’ shows the weighted rank, an indicator of relevancy to the prompt. Chunk A175X1 contains ‘financial~ statements’ glued together rather than separate; also its size is smaller than the top 2 chunks, making it more concentrated, and thus, potentially more interesting. By choosing specific chunks, you can dig deeper and retrieve the full content in the **detailed view** (not shown here). The **context view** shows the chunks with tags, titles, tables, images, and agents attached to them (not shown here) while the **embeddings view** in Figure 4 shows related keywords to use in subsequent prompts. The left column in that view shows the **PMI**, an indicator of relevancy. Note that ‘reconciliation’ is detected as a corpus word related to the prompt token ‘financial’. For the **agents view**, see Figure 2 in section 3.4.

```
-----
Prompt: financial statements 2024
Cleaned: ['2024', 'financial', 'statements']
-----
Most relevant text entities:
```

ID	wRank	size	PDF	pn	ID_Tokens
B2X0	5	6836	0	1	{'statements': 12, 'financial': 12, 'statements~financial': 1, '2024': 1}
A154X1	7	4083	1	158	{'statements': 4, 'financial': 4, 'financial~statements': 4, '2024': 1}
A175X1	7	2400	1	179	{'statements': 1, 'financial': 2, 'financial~statements': 1, '2024': 1}
A179X1	7	2687	1	183	{'statements': 1, 'financial': 1, 'financial~statements': 1, '2024': 1}
A182X1	7	1722	1	186	{'statements': 1, 'financial': 1, 'financial~statements': 1, '2024': 1}
A186X1	7	2319	1	190	{'statements': 1, 'financial': 1, 'financial~statements': 1, '2024': 2}
B2250X2	7	1992	2	85	{'statements': 1, 'financial': 1, 'financial~statements': 1, '2024': 3}
B2749X2	7	2662	2	127	{'statements': 2, 'financial': 4, 'financial~statements': 2, '2024': 1}

Figure 3: Sample prompt results, summary view

By allowing the user to choose frontend parameters in real-time, the system can collect the most popular parameter sets to automatically build a selection of default or template parameter sets. This is known as **self-tuning**. Finally, in the next version, it will be possible to give a higher weight to the first token in the prompt, and to favor the most recent results when timestamps are included in the corpus. Also, in the next version, the user will have the ability to browse next and previous chunks, starting at a chunk displayed in prompt results: this feature is easy to implement since chunks are sequentially indexed during the crawl. In short, the UI is

an **LLM browser** of its own, rather than a search box. Its scoring engine is described as the new **PageRank** for LLMs.

```
Top related tokens (via embeddings):
0.30 ('financial', 'non-gaap-gaap')
0.30 ('financial', 'gaap-financial')
0.30 ('financial', 'non-gaap-gaap-financial')
0.30 ('financial', 'gaap-financial-measures')
0.30 ('financial', 'reconciliation')
0.29 ('financial', 'financial-measures')
```

Figure 4: Sample prompt results, embeddings view

4.6 Testing and Evaluation

Evaluating LLMs is tricky for a number of reasons. Standard **evaluation** and benchmarking metrics fail to assess important qualities, such as **exhaustivity**, conciseness, depth and structuredness of prompt results. Part of the problem is because standard LLMs are trained to predict missing or next tokens, an important task in word guessing and earlier versions such as Bert, but now largely irrelevant to the problems it aims to solve.

It is compounded by the fact that wordy English prose is valued by many users, over concise, yet deep and well structured answers. Also, professional users and laymen will rate the results to a same prompt very differently, though a well designed LLM could deliver different results depending on the type of user. In our case, the user can choose agents or categories on the command menu when doing a search, and even **negative keywords**, to get more relevant results. Finally, our LLM shows **relevancy scores** attached to each item in prompt results. None of these features are taken into account in standard evaluation metrics.

To compare two different versions of our LLM, to check whether upgrades or different parameters actually lead to better performance, you want to separately test the different components impacted before a new release. More specifically:

- **Chunking.** Look for **chunks** that are too large or too small. Check the variance in chunk size, and also for the distribution of sub-chunk sizes within each parent chunk.
- **Stemming.** Try with and without **stemming** on the backend, with the **backend** parameter `use_stem` set to True or False. In the `hash_unstem` table, look for entries with 3 or more tokens attached to a stem, especially for stems with high occurrence; see whether they should be broken down into smaller pieces to make the stemming less aggressive and less ambiguous. Test with the **frontend** parameter `use_stem` set to True or False.
- **Duplicates.** Eliminate **duplicate chunks**, unless business reasons tell otherwise: if the same parent chunk is in two different PDFs, it might make sense to keep both copies. For sub-chunks, do not dedupe.
- **Distillation.** Done on the frontend. Test with and without **distillation**, by setting the frontend parameter `distill` to True or False.
- **Scoring.** Look at the **rank** and **score** attached to each chunk returned in the prompt results, for test prompts. Play with the related frontend parameters to **fine-tune** the scores. Are **g-tokens** and **s-tokens** given the proper weights? Are g-tokens correctly detected? Are **multi-tokens** favored over **single tokens**? Is token multiplicity and specialized tokens given the proper weights? Multiplicities are found in the **nested hash** `ID.hash[ID]`, where `ID` identifies a chunk, and `ID.hash[ID]` is the daughter hash storing its multi-tokens (the keys) that are also in the prompt with their multiplicities (the values).
- **Context.** Are there many chunks missing **contextual** fields, such as title or agent? Are tables correctly detected and indexed with the **multi-index** system? Are **agents** correctly assigned to chunks?
- **Speed.** Large chunks can significantly increase the size of backend tables `hash_pairs` and `embeddings`. These two tables also take more time to build, yet are not critical components of the system. I use them to suggest related prompts. Removing duplicate chunks requires `entity_list` which may take a lot of memory. Do we want to load all the chunks in memory (faster), especially the parent chunks which can be built using the sub-chunks? It is OK to have slower **backend** processing, as long as **frontend** is very fast. To improve speed and reduce memory usage, use lower values for the backend parameters `max_multitoken` and `maxDist`. Also, you can set `create_hpairs` to False: then no embedding is generated.

You may also try a different **PMI** function. PMIs are used in our **variable length embeddings**, as an alternative to **cosine similarity** and **dot product**. Also, you want to test with and without building the backend tables. Not building them means that they are pre-loaded in memory with `build_tables` set to False. You may want to verify that this step (reading the tables from text files and retrieving the correct architecture) works correctly.

As for actual evaluation, a possible test consists of extracting various pieces of text from the PDFs and check if they are retrieved depending on the prompts. In this case, the prompts are shorter, scrambled, cleaned versions of the original extract, possibly using different words. For instance, in the future, we will have a synonyms dictionary matching (say) ‘forecast’ to ‘projection’ and ‘prediction’, in case the word ‘forecast’ is in a prompt but not in the corpus. And even if in the corpus, to retrieve text containing not just ‘forecast’, but also its synonyms. Same with ‘gen ai’, ‘genai’ and ‘generative ai’.

5 Python code

The Python code is back-compatible with the previous version published in [1]. It still relies on fast **in-memory databases** stored as **nested hashes**, for backend tables. However, it has been significantly optimized for speed without increasing memory usage or table sizes, especially to deal with very large chunks. Finally, it has been significantly simplified, with some parameters no longer needed, and new ones added.

Now, you can read the backend tables from text files rather than creating them each time. Also, a number of functions have been moved to an internal library (see section 5.2) to make the code more readable. It also features new NLP functions related to **stemming** and **unstemming**.

5.1 Main program

The most recent version of the code is on a shared Google Drive, [here](#). It imports the internal library listed in section 5.2. The prompts are read from a text files available on GitHub, [here](#). A number of components need improvements, in particular the scoring engine, the stemmer, and g.tokens (gword in the code). Some need to be fully built, such as showing tables and images in the prompt results, overlapping chunks, more or better contextual fields, improved table detection, and showing absolute rather than relative scores in prompt results (the latter are just chunk ranks).

```

1  # xllm-enterprise-nvidia-dev.py
2  # text entity / sub-entity have same meaning as chunk / sub-chunk
3
4  import xllm_enterprise_nvidia_util as exllm
5
6  #--- Backend: create backend tables based on crawled corpus
7
8  # hash_pairs, ctokens are intermediate tables used when building backends
9  # no need to save them, or to load them when build_tables = False
10 # hash_context4 now contains IDs (short) rather than associated text (long)
11
12 tableNames = (
13     'dictionary', # multitokens (key = multitoken)
14     'hash_pairs', # multitoken associations (key = pairs of multitokens)
15     'ctokens',    # not adjacent pairs in hash_pairs (key = pairs of multitokens)
16     'hash_context1', # categories (key = multitoken)
17     'hash_context2', # tags (key = multitoken)
18     'hash_context3', # titles (key = multitoken)
19     'hash_context4', # chunk IDs attached to multitoken (key = multitoken)
20     'hash_context5', # meta (key = multitoken)
21     'hash_ID',      # text entity ID table (key = multitoken, value is list of IDs)
22     'hash_agents',  # agents (key = multitoken)
23     'ID_to_content', # full content attached to text entity ID (key = text entity ID)
24     'ID_to_agents',  # map text entity ID to agents list (key = text entity ID)
25     'ID_size',       # content size (key = text entity ID)
26     'ID_to_index',   # map ID to multi-index (key = text entity ID)
27     'Index_to_IDs',  # subIDs attached to ID, with size
28     'stopwords',     # stopwords list
29     'hash_stem',     # stemmed words in dictionary
30     'hash_unstem',   # match prompt words to multitokens (key = stemmed word)
31     'embeddings',    # to show related keywords in prompt results (key = multitoken)
32     'sorted_ngrams', # to build embeddings (key = multitoken),
33     # 'ID_to_tags',  # map ID to tag list (key = text entity ID) --- to be added
34 )
35
36 backendTables = {}
37 for name in tableNames:
38     backendTables[name] = {}
39
40 stopwords = ('', '-', 'in', 'the', 'and', 'to', 'of', 'a', 'this', 'for', 'is', 'with', 'from',
41             'as', 'on', 'an', 'that', 'it', 'are', 'within', 'will', 'by', 'or', 'its', 'can',
42             'your', 'be', 'about', 'used', 'our', 'their', 'you', 'into', 'using', 'these',
43             'which', 'we', 'how', 'see', 'below', 'all', 'use', 'across', 'provide', 'provides',

```



```

44         'aims', 'one', '&', 'ensuring', 'crucial', 'at', 'various', 'through', 'find', 'ensure',
45         'more', 'another', 'but', 'should', 'considered', 'provided', 'must', 'whether',
46         'located', 'where', 'begins', 'any', 'what', 'some', 'under', 'does', 'belong',
47         'included', 'part', 'associated')
48 backendTables['stopwords'] = stopwords
49
50 # agent_map key is corpus word, value is agent (many-to-one)
51 agent_map = { 'stock': 'Stock',
52               'fiscal 2022': 'Year 2022',
53               'fiscal 2023': 'Year 2023',
54               'income': 'Income',
55               'tax': 'Tax',
56               'cash': 'Cash',
57               'products': 'Products',
58               'revenue': 'Revenue',
59               'directors': 'Directors',
60               'data': 'Data',
61               'equity': 'Equity',
62               'management': 'Management',
63               'common stock': 'Common Stock',
64               'securities': 'Securities',
65               'assets': 'Assests',
66               'financial statements': 'Financial Statements',
67               'restricted stock': 'Restricted Stock',
68               'accelerated computing': 'Accelerated Computing',
69               'data center': 'Data Center',
70               'non-gaap': 'Non-GAAP',
71             }
72
73
74 #--- Read repository and create all backend tables
75
76 def read_backend_tables(backendTables):
77
78     for tableName in backendTables:
79
80         print("reading %16s" %(tableName), end = " ")
81         filename = "backend_" + tableName + ".txt"
82         if tableName == 'stopwords':
83             backendTables[tableName] = exllm.read_list(filename)
84         elif tableName in ('dictionary', 'ID_size'):
85             backendTables[tableName] = exllm.read_pairs(filename)
86         elif tableName in ('hash_stem', ):
87             backendTables[tableName] = exllm.read_table(filename, format = "str")
88         elif tableName not in ('hash_pairs', 'ctokens'):
89             backendTables[tableName] = exllm.read_table(filename)
90         print("(size: %8d)" %(len(backendTables[tableName])))
91
92     return(backendTables)
93
94
95 def build_backend_tables(repository, backendTables, save=True):
96
97     backendParams = {
98         'max_multitoken': 4, # max. consecutive terms per multi-token for inclusion in dictionary
99         'maxDist' : 3,      # max. position delta between 2 multitokens to link them in hash_pairs
100        'create_hpairs': True, # required for embeddings, takes time and memory
101        'create_ctokens': False, # needs create_hpairs set to TRUE for activation
102        'use_stem': False, # backend stemming
103        'extraWeights' : # default weight is 1
104            {
105                'description': 0.0,
106                'category': 0.0, # 0.5
107                'tag_list': 0.0, # 0.5
108                'title': 0.0, # 0.5
109                'meta': 0.0 # 0.5
110            }
111    }
112
113    IN = open(repository, "r")
114    data = IN.read()
115    IN.close()
116
117    entities = data.split("\n")
118    count = 0
119    context_fields = ('category_text', 'tags_list_text', 'title_text')

```

```

120 ID_size = backendTables['ID_size']
121
122 # to avoid duplicate entities (takes space, better to remove them in the corpus)
123 entity_list = ()
124
125 for entity_raw in entities:
126
127     entity = entity_raw.split("~")
128     agent_list = ()
129
130     if len(entity) > 1 and entity[1] not in entity_list:
131
132         n1 = len(entities)
133         n2 = len(backendTables['dictionary'])
134         n3 = len(backendTables['hash_pairs'])
135         if count % 50 == 0:
136             print("Processing chunk %4d out of %5d [%6d %6d]" %(count,n1, n2, n3))
137         count += 1
138
139         # entity_list = (*entity_list, entity[1]) # avoid duplicates [use lots memory]
140         hash_chunk = eval(entity[1])
141
142         # entity_ID can be a string (multi-index)
143         entity_ID = entity[0]
144         hash_crawl = {}
145         hash_crawl['ID'] = entity_ID
146         ID_size[entity_ID] = len(entity[1])
147
148         for key in hash_chunk:
149
150             value = hash_chunk[key]
151             if key == 'category_text':
152                 hash_crawl['category'] = value
153             elif key == 'tags_list_text':
154                 hash_crawl['tag_list'] = exllm.clean_list(value)
155             elif key == 'title_text':
156                 hash_crawl['title'] = value
157             elif key == 'description_text':
158                 hash_crawl['description'] = value # do not build to save space
159             elif key == 'type':
160                 hash_crawl['meta'] = value
161             if key in context_fields or key in 'description_text':
162                 # remove <in 'description_text'>, except in demo
163                 for word in agent_map:
164                     agent = agent_map[word]
165                     if key == 'description_text':
166                         agent = agent.lower()
167                     if word in value.lower() and agent not in agent_list:
168                         agent_list = (*agent_list, agent)
169
170             # hash_crawl['full_content'] = hash_chunk
171             hash_crawl['mindex'] = (hash_chunk['doc_ID'], hash_chunk['pn'])
172             hash_crawl['agents'] = agent_list
173             exllm.update_dict(backendTables, hash_crawl, backendParams)
174
175 #-- Create embeddings
176
177 embeddings = {} # multitoken embeddings based on hash_pairs
178
179 hash_pairs = backendTables['hash_pairs']
180 dictionary = backendTables['dictionary']
181
182 for key in hash_pairs:
183     wordA = key[0]
184     wordB = key[1]
185     nA = dictionary[wordA]
186     nB = dictionary[wordB]
187     nAB = hash_pairs[key]
188     pmi = nAB/(nA*nB)**0.5 # try: nAB/(nA + nB - nAB)
189     # if nA + nB <= nAB:
190     #     print(key, nA, nB, nAB)
191     exllm.update_nestedHash(embeddings, wordA, wordB, pmi)
192     exllm.update_nestedHash(embeddings, wordB, wordA, pmi)
193
194 backendTables['embeddings'] = embeddings
195

```

```

196     #-- Create sorted n-grams
197
198     sorted_ngrams = {} # to match ngram prompts with embeddings entries
199
200     for word in dictionary:
201         tokens = word.split('~')
202         tokens.sort()
203         sorted_ngram = tokens[0]
204         for token in tokens[1:len(tokens)]:
205             sorted_ngram += "~" + token
206         exllm.update_listHash(sorted_ngrams, sorted_ngram, word)
207
208     backendTables['sorted_ngrams'] = sorted_ngrams
209
210     #-- Create stem table
211
212     (hash_stem, hash_unstem) = exllm.stem(dictionary)
213     for word in hash_stem:
214         # need to manually delete some entries [not done yet]
215         lead_word = hash_stem[word]
216         cnt1 = int(dictionary[word])
217         cnt2 = int(dictionary[lead_word])
218         print("%3d\t%2d\t%s\t%s" %(cnt1, cnt2, word, lead_word))
219     backendTables['hash_stem'] = hash_stem
220     backendTables['hash_unstem'] = hash_unstem
221
222     #-- save backend tables
223
224     if save:
225
226         # save backend tables
227         for tableName in backendTables:
228             if tableName not in ('hash_pairs', 'ctokens'):
229                 table = backendTables[tableName]
230                 print("Saving table %16s %8d elts %s" %(tableName, len(table), type(table)))
231                 OUT = open('backend_' + tableName + '.txt', "w")
232                 if type(table) is tuple:
233                     OUT.write(str(table))
234                 else:
235                     for key in table:
236                         value = table[key]
237                         OUT.write(str(key) + "\t" + str(value) + "\n")
238                 OUT.close()
239
240         # save backend parameters
241         OUT = open('backendParams.txt', "w")
242         OUT.write(str(backendParams))
243         OUT.close()
244
245     return(backendTables)
246
247
248     #--- Main for backend (reading or creating backend tables)
249
250     build_tables = True
251
252     if build_tables:
253         repository_file = "repository_Nvidia_X.txt"
254         build_backend_tables(repository_file, backendTables)
255         embeddings = backendTables['embeddings']
256         sorted_ngrams = backendTables['sorted_ngrams']
257     else:
258         read_backend_tables(backendTables)
259
260     dictionary = backendTables['dictionary']
261     DICT = open("dict.txt", "wt")
262
263     for key in dictionary:
264         count = dictionary[key]
265         ntk = 1 + key.count("~") # number of single tokens in multitokens
266         DICT.write(str(count) + "\t" + str(ntk) + "\t" + key + "\n")
267
268     DICT.close()
269
270
271     #--- Functions used to score results ---

```

```

272
273 def rank(hash):
274     # sort hash, then replace values with their rank
275
276     hash = dict(sorted(hash.items(), key=lambda item: item[1], reverse=True))
277     rank = 0
278     old_value = 999999999999
279
280     for key in hash:
281         value = hash[key]
282         if value < old_value:
283             rank += 1
284             hash[key] = rank
285             old_value = value
286     return(hash)
287
288
289 def rank_ID(ID_score):
290     # attach weighted relevancy rank to text entity ID, with respect to prompt
291
292     ID_score0 = {}
293     ID_score1 = {}
294     ID_score2 = {}
295     ID_score3 = {}
296
297     for ID in ID_score:
298         score = ID_score[ID]
299         ID_score0[ID] = score[0]
300         ID_score1[ID] = score[1]
301         ID_score2[ID] = score[2]
302         ID_score3[ID] = score[3]
303
304     ID_score0 = rank(ID_score0)
305     ID_score1 = rank(ID_score1)
306     ID_score2 = rank(ID_score2)
307     ID_score3 = rank(ID_score3)
308
309     ID_score_ranked = {}
310     for ID in ID_score:
311         weighted_rank = 2*ID_score0[ID] + ID_score1[ID] + ID_score2[ID] + ID_score3[ID]
312         ID_score_ranked[ID] = weighted_rank
313     ID_score_ranked = dict(sorted(ID_score_ranked.items(), key=lambda item: item[1]))
314     return(ID_score_ranked)
315
316
317 #--- Main for frontend (prompt processing prompts)
318
319 print("\n")
320 input_ = " "
321 saved_query = ""
322 get_bin = lambda x, n: format(x, 'b').zfill(n)
323 frontendParams = exllm.default_frontendParams()
324 use_stem = frontendParams['use_stem']
325 beta = frontendParams['beta']
326 ID_to_content = backendTables['ID_to_content']
327
328
329 #--- Main: Read sample prompts ---
330
331 from nltk.stem import PorterStemmer
332 from collections import defaultdict
333 stemmer = PorterStemmer()
334
335 sorted_ngrams = backendTables['sorted_ngrams']
336 embeddings = backendTables['embeddings']
337 ID_to_content = backendTables['ID_to_content']
338 hash_unstem = backendTables['hash_unstem']
339
340 IN = open("enterprise_nvidia_prompts.txt", "r")
341 prompts = IN.read()
342 prompts = prompts.split("\n")
343
344 # --- Main: Look over all prompts ---
345
346 for query in prompts:
347

```

```

348 query = query.split("|")[0]
349 print("\n-----")
350 print("Prompt: ", query)
351 query = query.replace('?', ' ').replace('(', ' ').replace(')', ' ').replace('.', ' ')
352 query = query.replace('"', '').replace("\s", '')
353 query = query.split(' ')
354 for k in range(len(query)):
355     query[k] = query[k].lower() # need to eliminate this down the line
356 new_query = []
357 neg_query = [] # keywords to exclude
358 altTokens = ()
359
360 for k in range(len(query)):
361     token = query[k]
362     if token[0] == '!':
363         token = token[1:len(token)]
364         neg_query.append(token)
365     if use_stem:
366         tstem = stemmer.stem(token)
367         if tstem in hash_unstem:
368             tlist = hash_unstem[tstem]
369             for altToken in tlist:
370                 if token != altToken and altToken not in altTokens :
371                     altTokens = (*altTokens, altToken)
372     if token in dictionary:
373         new_query.append(token)
374
375 q_altTokens = ()
376 for altToken in altTokens:
377     if altToken not in new_query:
378         new_query.append(altToken)
379     q_altTokens = (*q_altTokens, altToken)
380
381 query = new_query.copy()
382 query.sort()
383 print("Cleaned:", query)
384 print("-----")
385
386 q_embeddings = {}
387 q_dictionary = {}
388
389 # --- build q_dictionary and q_embeddings based on prompt tokens ---
390
391 for k in range(1, 2*len(query)):
392
393     binary = get_bin(k, len(query))
394     sorted_word = ""
395     for k in range(0, len(binary)):
396         if binary[k] == '1':
397             if sorted_word == "":
398                 sorted_word = query[k]
399             else:
400                 sorted_word += "~" + query[k]
401
402     if sorted_word in sorted_ngrams:
403         ngrams = sorted_ngrams[sorted_word]
404         for word in ngrams:
405             if word in dictionary:
406                 q_dictionary[word] = dictionary[word]
407                 embedding = exllm.get_value(word, embeddings)
408                 exllm.add_embedding(q_embeddings, word, embedding)
409
410 # deal with prompt multitokens, if there are any [need to add multitoken stemming]
411 for word in query:
412     if '~' in word and word in dictionary and word not in q_dictionary:
413         q_dictionary[word] = dictionary[word]
414         embedding = exllm.get_value(word, embeddings)
415         exllm.add_embedding(q_embeddings, word, embedding)
416
417 # --- Scoring and selecting what to show in prompt results ---
418
419 if frontendParams['distill']:
420     # how is this working with negative keywords?
421     exllm.distill_frontendTables(q_dictionary, q_embeddings, frontendParams)
422 hash_ID = backendTables['hash_ID']
423 ID_hash = {} # local, transposed of hash_ID; key = ID; value = multitoken list

```

```

424
425 for word in q_dictionary:
426     for ID in hash_ID[word]:
427         local_hash = hash_ID[word]
428         if word not in neg_query:
429             exllm.update_nestedHash(ID_hash, ID, word, local_hash[ID])
430         gword = "__" + word # graph multitoken
431         if gword in hash_ID and word not in neg_query:
432             for ID in hash_ID[gword]:
433                 exllm.update_nestedHash(ID_hash, ID, gword, 1)
434
435 ID_score = {}
436 for ID in ID_hash:
437     # score[0] is inverse weighted count
438     # score[1] is raw number of tokens found
439     score = [0, 0] # based on tokens present in the entire text entity
440     gscore = [0, 0] # based on tokens present in graph (context elements)
441     for token in ID_hash[ID]:
442         if token in dictionary:
443             score[0] += 1/(q_dictionary[token]**beta)
444             score[1] += 1
445         else:
446             # token must start with "__" (it's a graph token)
447             token = token[2:len(token)]
448             gscore[0] += 1/(q_dictionary[token]**beta)
449             gscore[1] += 1
450     ID_score[ID] = [score[0], score[1], gscore[0], gscore[1]]
451
452 # --- Print results ---
453
454 ID_score_ranked = rank_ID(ID_score)
455 nresults = frontendParams['nresults']
456
457 print("Most relevant chunks with multitokens, doc_ID, pn, rank, size:\n")
458 # also show absolute score as opposed to rank
459 n_ID = 0
460 print("\n      ID wRank size PDF pn ID_Tokens")
461 ID_index = backendTables['ID_to_index']
462 ID_size = backendTables['ID_size']
463 for ID in ID_score_ranked:
464     if n_ID < nresults:
465         # content of text entity ID not shown, stored in ID_to_content[ID]
466         # add tags
467         minindex = ID_index[ID] # multi-index
468         doc_ID = minindex[0] # PDF number
469         pn = minindex[1] # page number within PDF
470         rankx = ID_score_ranked[ID] # need to also create ID_score[ID] (absolute score)
471         size = ID_size[ID]
472         print(" %8s %3d %6d %3d %4d %s"
473               %(ID, rankx, size, doc_ID, pn, ID_hash[ID]))
474         n_ID += 1
475
476 print("Most relevant chunks with agents:\n")
477 n_ID = 0
478 print("\n      ID      ID_agents")
479 ID_to_agents = backendTables['ID_to_agents']
480 for ID in ID_score_ranked:
481     if n_ID < nresults:
482         agents = exllm.get_value(ID, ID_to_agents)
483         if len(agents) > 0:
484             print(" %8s %s" %(ID, agents))
485         n_ID += 1
486
487 print("\nToken count (via dictionary):\n")
488 for key in q_dictionary:
489     print(" %4d %s" %(q_dictionary[key], key))
490
491 print("\nTop related tokens (via embeddings):\n")
492 q_embeddings = dict(sorted(q_embeddings.items(), key=lambda item: item[1], reverse=True))
493 n_words = 0
494 for word in q_embeddings:
495     pmi = q_embeddings[word]
496     if n_words < 10:
497         print(" %5.2f %s" %(pmi, word))
498     n_words += 1
499

```

```

500 full_content = False
501 if full_content:
502     print("\nFull content sorted by relevancy\n")
503     # ID starts with A for big chunks, with B for small chunks
504     n_ID = 0
505     for ID in ID_score_ranked:
506         content = ID_to_content[ID]
507         rankx = ID_score_ranked[ID]
508         size = ID_size[ID]
509         if n_ID < nresults and size < 60000:
510             print("%8s %3d %5d %s %s\n" % (ID, rankx, size, ID_hash[ID], content))
511             n_ID += 1

```

5.2 Internal library

This library is on Google Drive, [here](#). The backend tables are also in the same GitHub folder. In this version, I removed the code for command-line processing (real-time fine-tuning and so on) as prompts are read from a text file. However, it is included in the previous version featured in Part 1 in [1] and the corresponding features will be accessible from the UI in the new Web API under construction.

```

1  # xllm_enterprise_nvidia_util.py
2
3  #--- Stemming
4
5  def stem(dictionary):
6
7      from nltk.stem import PorterStemmer
8      # from nltk.stem import WordNetLemmatizer
9      stemmer = PorterStemmer()
10     # lemmatizer = WordNetLemmatizer()
11
12     hash_unstem = {}
13     for word in dictionary:
14         if word.count("'") == 0:
15             key = stemmer.stem(word)
16             # key = lemmatizer.lemmatize(word)
17             hash_unstem = update_listHash(hash_unstem, key, word)
18
19     hash_stem = {}
20     for key in hash_unstem:
21         list = hash_unstem[key]
22         if len(list) > 1:
23             # a few stems with 3+ words need breaking down [not done yet]
24             # print(key, hash_unstem[key])
25             max_cnt = 0
26             for word in list:
27                 cnt = dictionary[word]
28                 if cnt > max_cnt:
29                     max_cnt = cnt
30                     lead_word = word
31             for word in list:
32                 if word != lead_word:
33                     hash_stem[word] = lead_word
34
35     hash_stem = dict(sorted(hash_stem.items()))
36     return(hash_stem, hash_unstem)
37
38
39 #--- Read backend-tables
40
41 def get_data(filename, path):
42     if 'http' in path:
43         response = requests.get(path + filename)
44         data = (response.text).replace('\r', '').split("\n")
45     else:
46         file = open(filename, "r")
47         data = [line.rstrip() for line in file.readlines()]
48         file.close()
49     return(data)
50
51
52 def read_table(filename, format = "float", path = ''):
53     table = {}

```

```

54 data = get_data(filename, path)
55 for line in data:
56     line = line.split('\t')
57     if len(line) > 1:
58         value = line[1]
59         if format == 'str':
60             table[line[0]] = line[1]
61         else:
62             table[line[0]] = eval(line[1])
63 return(table)
64
65
66 def read_list(filename, path = ''):
67     data = get_data(filename, path)
68     stopwords = eval(data[0])
69     return(stopwords)
70
71
72 def read_pairs(filename, path = ''):
73     dictionary = {}
74     data = get_data(filename, path)
75     for line in data:
76         line = line.split('\t')
77         if len(line) > 1:
78             if type(line[1]) is not float:
79                 dictionary[line[0]] = float(line[1])
80             else:
81                 dictionary[line[0]] = line[1]
82     return(dictionary)
83
84
85 #--- Hash functions
86
87 def update_listHash(hash, key, value):
88
89     if key in hash:
90         list = hash[key]
91         if value not in list:
92             list = (*list, value)
93     else:
94         list = (value,)
95     hash[key] = list
96     return(hash)
97
98
99 def update_hash(hash, key, count=1):
100
101     if key in hash:
102         hash[key] += count
103     else:
104         hash[key] = count
105     return(hash)
106
107
108 def update_nestedHash(hash, key, value, count=1):
109
110     # 'key' is a word here, value is tuple or single value
111     if key in hash:
112         local_hash = hash[key]
113     else:
114         local_hash = {}
115     if type(value) is not tuple:
116         value = (value,)
117     for item in value:
118         if item in local_hash:
119             local_hash[item] += count
120         else:
121             local_hash[item] = count
122     hash[key] = local_hash
123     return(hash)
124
125
126 def get_value(key, hash):
127     if key in hash:
128         value = hash[key]
129     else:

```



```

130     value = ''
131     return(value)
132
133
134 #--- Build back-end tables
135
136 def update_tables(backendTables, word, hash_crawl, backendParams):
137
138     category = get_value('category', hash_crawl)
139     tag_list = get_value('tag_list', hash_crawl)
140     title = get_value('title', hash_crawl)
141     description = get_value('description', hash_crawl) #
142     meta = get_value('meta', hash_crawl)
143     ID = get_value('ID', hash_crawl)
144     agents = get_value('agents', hash_crawl)
145     full_content = get_value('full_content', hash_crawl) #
146     mindex = get_value('mindex', hash_crawl)
147
148     ID_size = backendTables['ID_size']
149
150     extraWeights = backendParams['extraWeights']
151     word = word.lower() # add stemming
152     weight = 1.0
153     flag = ''
154     if word in category:
155         weight += extraWeights['category']
156         flag = '___'
157     if word in tag_list:
158         weight += extraWeights['tag_list']
159         flag = '___'
160     if word in title:
161         weight += extraWeights['title']
162         flag = '___'
163     if word in meta:
164         weight += extraWeights['meta']
165         flag = '___'
166
167     if flag != '':
168         gword = flag + word
169         update_nestedHash(backendTables['hash_ID'], gword, ID)
170
171     update_hash(backendTables['dictionary'], word, weight)
172     update_nestedHash(backendTables['hash_context1'], word, category)
173     update_nestedHash(backendTables['hash_context2'], word, tag_list)
174     update_nestedHash(backendTables['hash_context3'], word, title)
175     update_nestedHash(backendTables['hash_context4'], word, ID) # used to be 'description'
176     update_nestedHash(backendTables['hash_context5'], word, meta)
177     update_nestedHash(backendTables['hash_ID'], word, ID)
178     update_nestedHash(backendTables['hash_agents'], word, agents)
179     # update_nestedHash(backendTables['full_content'], word, full_content) # takes space, don't
        nuuild?
180
181     if ID not in backendTables['ID_to_content']:
182         for agent in agents:
183             # new format: listHash; old format: nestedHash
184             # update_nestedHash(backendTables['ID_to_agents'], ID, agent)
185             update_listHash(backendTables['ID_to_agents'], ID, agent)
186             update_hash(backendTables['ID_to_content'], ID, full_content)
187             update_hash(backendTables['ID_to_index'], ID, mindex)
188             update_nestedHash(backendTables['Index_to_IDs'], mindex, ID, ID_size[ID])
189
190     return(backendTables)
191
192
193 def clean_list(value):
194
195     # change string "[ 'a', 'b', ...]" to ('a', 'b', ...)
196     value = value.replace("[", "").replace("]", "")
197     aux = value.split("~")
198     value_list = ()
199     for val in aux:
200         val = val.replace("'", "").replace('"', "").rstrip()
201         if val != '':
202             value_list = (*value_list, val)
203     return(value_list)
204

```

```

205
206 def get_key_value_pairs(entity):
207
208     # extract key-value pairs from 'entity' (a string)
209     entity = entity[1].replace("}", "{, '")
210     flag = False
211     entity2 = ""
212
213     for idx in range(len(entity)):
214         if entity[idx] == '[':
215             flag = True
216         elif entity[idx] == ']':
217             flag = False
218         if flag and entity[idx] == ",":
219             entity2 += "~"
220         else:
221             entity2 += entity[idx]
222
223     entity = entity2
224     key_value_pairs = entity.split(", ")
225     return(key_value_pairs)
226
227
228 def update_dict(backendTables, hash_crawl, backendParams):
229
230     max_multitoken = backendParams['max_multitoken']
231     maxDist = backendParams['maxDist']
232     create_hpairs = backendParams['create_hpairs']
233     create_ctokens = backendParams['create_ctokens']
234     use_stem = backendParams['use_stem']
235
236
237     category = get_value('category', hash_crawl)
238     tag_list = get_value('tag_list', hash_crawl)
239     title = get_value('title', hash_crawl)
240     description = get_value('description', hash_crawl)
241     meta = get_value('meta', hash_crawl)
242
243     text = category + "." + str(tag_list) + "." + title + "." + description + "."
244     text = text.replace('/', " ").replace('(', ' ').replace(')', ' ').replace('?', '')
245     text = text.replace(" ", "").replace("'", "").replace('\\n', '').replace('!', ' ')
246     text = text.replace("\\s", '').replace("\\t", '').replace(", ", " ").replace(":", " ")
247     text = text.replace(";", " ").replace("|", " ").replace("--", " ").replace(" ", " ").lower()
248     sentence_separators = ('.',)
249     for sep in sentence_separators:
250         text = text.replace(sep, '~')
251     text = text.split('~')
252
253     hash_pairs = backendTables['hash_pairs']
254     ctokens = backendTables['ctokens']
255     hash_stem = backendTables['hash_stem']
256     stopwords = backendTables['stopwords']
257     buffer2 = []
258
259     for sentence in text:
260
261         words = sentence.split(" ")
262         offset = 0
263         buffer = []
264
265         for word in words:
266
267             if use_stem and word in hash_stem:
268                 word = hash_stem[word]
269
270             if word not in stopwords:
271                 # word is single token
272                 buffer.append(word)
273                 buffer2.append(word)
274                 update_tables(backendTables, word, hash_crawl, backendParams)
275
276             for k in range(1, max_multitoken):
277                 if offset > 0:
278                     # word is now multi-token with k+1 tokens
279                     word = buffer[offset-k] + "~" + word
280                     buffer2.append(word)

```

```

281         update_tables(backendTables, word, hash_crawl, backendParams)
282
283         offset += 1
284
285     if create_hpairs:
286
287         for k in range(len(buffer2)):
288
289             wordA = buffer2[k]
290             lbound = max(0, k-maxDist) # try bigger value for maxDist
291             ubound = min(len(buffer2), k+maxDist+1)
292
293             for l in range(lbound, ubound):
294                 wordB = buffer2[l]
295                 key = (wordA, wordB)
296                 if wordA < wordB:
297                     hash_pairs = update_hash(hash_pairs, key)
298                     if create_ctokens and k != 1:
299                         ctokens = update_hash(ctokens, key)
300
301         return(backendTables)
302
303
304 def default_frontendParams():
305
306     frontendParams = {
307         'distill': False,
308         'maxTokenCount': 1000, # ignore generic tokens if large enough
309         'beta': 1.0, # used in text entity relevancy score, try 0.5
310         'nresults': 20, # max number of chunks to show in results
311         'use_stem': False,
312     }
313     return(frontendParams)
314
315
316 def distill_frontendTables(q_dictionary, q_embeddings, frontendParams):
317     # purge q_dictionary then q_embeddings (frontend tables)
318
319     maxTokenCount = frontendParams['maxTokenCount']
320     local_hash = {}
321     for key in q_dictionary:
322         if q_dictionary[key] > maxTokenCount:
323             local_hash[key] = 1
324     for keyA in q_dictionary:
325         for keyB in q_dictionary:
326             nA = q_dictionary[keyA]
327             nB = q_dictionary[keyB]
328             if keyA != keyB:
329                 if (keyA in keyB and nA == nB) or (keyA in keyB.split('~')):
330                     local_hash[keyA] = 1
331     for key in local_hash:
332         del q_dictionary[key]
333
334     local_hash = {}
335     for key in q_embeddings:
336         if key[0] not in q_dictionary:
337             local_hash[key] = 1
338     for key in local_hash:
339         del q_embeddings[key]
340
341     return(q_dictionary, q_embeddings)
342
343
344 def add_embedding(q_embeddings, word, embedding):
345     for token in embedding:
346         pmi = embedding[token]
347         q_embeddings[(word, token)] = float(pmi)
348     return(q_embeddings)

```

References

- [1] Vincent Granville. *Building Disruptive AI & LLM Technology from Scratch*. MLTechniques.com, 2024. [\[Link\]](#). 1, 2, 3, 9, 12, 14, 16, 23

- [2] Vincent Granville. *State of the Art in GenAI & LLMs – Creative Projects, with Solutions*. MLTechniques.com, 2024. [[Link](#)]. [1](#), [2](#)

Index

- n*-grams
 - sorted, 12
- agent, 2, 11, 12, 15
 - action agent, 11
 - hard agent, 12
 - retrieval agent, 11, 12
 - soft agent, 12
- backend, 9, 15
- c_token, 9
- cataloging, 1
- chunking, 2, 12, 15
 - backen tables, 9
 - duplicates, 15
 - overlap, 3
 - parent chunk, 11
 - sub-chunk, 2, 11
- clustering
 - text clustering, 1
- context, 15
 - context features, 10
 - contextual window, 10
- cosine similarity, 15
- distillation, 10, 12
 - frontend, 15
- dot product, 15
- embedding, 12
 - variable length, 9
- embeddings
 - variable length, 15
- evaluation, 15
- exhaustivity, 15
- explainable AI, 9, 13
- fine-tuning, 15
 - real-time, 12, 14
 - self-tuning, 14
- frontend, 9, 12, 15
- g_token, 9
- hallucination, 12
- hash table
 - in-memory, 16
 - nested hash, 10, 12, 15, 16
- in-memory processing, 16
- index
 - multi-index, 2, 9, 12, 15
- knowledge graph, 9
- LLM
 - debugging, 13
 - LLM browser, 15
 - LLM router, 1
 - sub-LLM, 1
- xLLM, 9
 - enterprise version, 9
- multimodal, 2
- negative keyword, 12, 14, 15
- NLTK, 13
- PMI, 9, 14, 15
- prompt engineering, 9
- PymuPDF, 3
- relevancy, 15
 - LLM PageRank, 15
 - rank, 12, 15
 - score, 12, 14, 15
 - score vector, 12
 - scoring engine, 10
 - sub-score, 12
- reproducibility, 9
- s_token, 9
- stemming, 15, 16
 - backend, 11
 - frontend, 13
 - unstemming, 13
- synthetic data, 1
- token
 - g_token, 15
 - multi-token, 9, 14, 15
 - s_token, 15
 - single token, 15
- transformer, 13
- unstemming, 12, 16
- view
 - agents, 12, 14
 - context, 14
 - detailed, 14
 - embeddings, 14
 - summary, 14