Numbers Matter! Bringing Quantity-awareness to Retrieval Systems

Anonymous ACL submission

Abstract

Quantitative information plays a crucial role in understanding and interpreting content of documents. Many user queries contain quantities and cannot be resolved without understanding their semantics, e.g., "car that costs less than \$10k". Yet, modern search engines apply the same ranking mechanisms for both words and quantities, overlooking magnitude and unit information. In this paper, we introduce two quantity-aware ranking techniques designed to rank both the quantity and textual content either jointly or independently. These techniques incorporate quantity information in available retrieval systems and can address queries with numerical conditions equal, greater than, and less than. To evaluate the effectiveness of our proposed models, we introduce two novel quantity-aware benchmark datasets in the domains of finance and medicine and compare our method against various lexical and neural models. The code and data are available under https://github.com/filled_in_later.

1 Introduction

011

013

017

019

021

033

037

041

Despite advances in semantic search, and sophisticated neural network architectures, handling quantitative information in text remains challenging. Specifically with quantity-centric queries, in which the query contains a quantity and a numerical condition, e.g., "BMW with more than 530hp". The reason for this is that systems are not aware of numbers and their semantics, such as proximity, in particular in combination with units. Numbers and units are treated in the same way as any other text token that is subject to subsequent processing, e.g., indexing or embedding. What complicates treating numbers and units in a proper way is that these objects can also have different surface forms (e.g., 6k vs 6,000 and mph vs miles per hour) and require standardization (Weikum, 2020). While there are approaches that specifically focus on numbers in text, e.g., extracting quantities for entities (Ho et al.,

2019; Li et al., 2021), linking quantities in tables (Ibrahim et al., 2019), or numerical reasoning (Ran et al., 2019), they are tailored for specific tasks and not semantic search in general. This applies to neural models supporting Information Retrieval (IR), which are trained on general-purpose data without the focus on quantity semantics. Language Models (LM), forming the basis for neural models, exhibit a limited understanding of number scales and proximity (Wallace et al., 2019). Despite recent work on numerical language models (Spokoyny et al., 2022; Jin et al., 2021), these architectures are very specific and require changes in the architecture of popularly used language models in IR, which indicates an expensive pre-training. Moreover, lack of accessible quantity-centric benchmarks for training or comparing systems exacerbates the issue.

042

043

044

047

048

053

054

056

060

061

062

063

064

065

066

067

068

069

070

071

072

073

074

076

077

078

079

081

In this paper, we present two strategies to enhance the quantity understanding of current IR systems. We aim for a general-purpose model that is not specific to quantity ranking but is also capable of textual ranking. The two approaches differ in their integration of quantity ranking with textual ranking. The first employs a disjoint combination, while the second focuses on the joint ranking of quantities within the context of textual content. The disjoint approach is an unsupervised and heuristic model utilizing an index structure, compatible with various lexical and semantic IR systems. Due to the independence assumption, the connection between quantities and surrounding text is somewhat lost. Therefore, for joint ranking, we aim to learn quantity-aware document and query representations through task-specific fine-tuning of neural IR models. Additionally, we introduce two novel benchmark datasets for quantity-centric ranking, specifically focusing on queries involving numerical conditions, in the domains of finance and medicine. We evaluate the performance of our systems against various lexical and neural models and show significant improvements over the baselines.

2 Related Work

086

880

089

094

098

100

101

104

105

106

131

Related work for quantity-centric search is limited. (Ho et al., 2020, 2019) focuses on quantity search for named entities, using a deep neural network for extracting quantity-centric tuples from text and query and matching based on context similarity. Their pipeline involves semantic role labeling and named entity extraction, both resource-intensive and reliant on sparsely available annotated data for quantities. Further, focusing on named entities limits the applicability to real-world scenarios.

QFinder (Almasian et al., 2022) integrates numerical and lexical indexes to enhance numerical understanding in a lexical IR system. Our disjoint model utilizes QFinder's heuristic ranking function, but we extend their approach to include neural models and go beyond the limited query language, allowing users to provide queries in plain text. MQSearch (Maiya et al., 2015) extracts quantities with a set of regular expressions to create a rule-based system for finding documents containing certain keywords and ranges of values. Loosely related to IR, (Rybinski et al., 2023) and (Li et al., 2021) perform numerical summarization on unstructured text in form of plots and graphs.

108 In the area of databases, there has been some work focusing on building numerical indices for queries 109 that contain numerical restrictions (Maiya et al., 110 2015; Fontoura et al., 2007; Agrawal and Srikant, 111 2003). However, the main focus of such systems 112 is the efficiency of the index structure and filtering 113 out irrelevant numbers from the results with hard 114 constraints rather than ranking. 115

Unlike quantity-aware IR, investigating numeracy 116 in LMs is well-established. (Wallace et al., 2019) 117 is among the first to highlight the limitations of 118 embedding models when handling numbers. Subse-119 quent efforts have led to dedicated embeddings and LMs for understanding scales, basic arithmetic, and 121 numerical common sense knowledge (Spokoyny 122 et al., 2022; Jin et al., 2021; Thawani et al., 2021; 123 Sundararaman et al., 2020; Jiang et al., 2020; 124 Nogueira et al., 2021; Spithourakis and Riedel, 125 2018). These models are specific to numeracy and 126 not IR in general. While using them can enhance 127 performance, we focus on improving quantity un-128 129 derstanding in current IR models without architecture change or training a LM from scratch. 130

3 Quantity-aware Model

A quantity-centric query contains a numerical condition, a value, and a unit, e.g., "iPhone XS with price under \$1500 ". Queries like "What is the price of iPhone XS?" are not considered quantitycentric as they don't require an understanding of scales and units. In the following, we assume a document collection where each document consists of a sequence of sentences. Following previous work (Ho et al., 2019; Almasian et al., 2022), we focus on sentences as retrieval units. A sentence $s_i := (T_i, Q_i)$ is a sequence of tokens $T_i = (t_1, ..., t_l)$ and quantities $Q_i = (q_1, ..., q_k)$, where a quantity $q_i = (u_i, v_i)$ is a tuple of a unit u_i and a value v_i . A quantity query is denoted by $X = (T_x, c, q_x)$, where $T_x = (t_{x_1}, .., t_{x_n})$ are the search terms related to the query quantity $q_x = (u_x, v_x)$. $c \in \{=, <, >\}$ represents a numerical condition, defining equal, less than, and greater than conditions. Less than and greater than indicate open bounds with values strictly less or greater than the query value. The equal condition pertains to values strictly equal to the query value. The relevance, $r(s_i|X)$, of sentence s_i to query X is denoted in Eq 1. The similarity function sim_c is dependent on the query condition c, where τ is a generic function that maps a query and document to their representations. Here, we explore different ways to define τ , which can be an embedding vector or a heuristic scoring function.

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

149

150

151

152

153

154

155

156

157

158

160

162

163

164

165

166

167

169

170

171

172

173

174

175

176

177

178

179

180

183

$$r(s_i|X) \sim sim_c(\tau(T_x, q_x), \tau(T_i, Q_i))$$
(1)

We begin with a disjoint quantity-ranking method. Leveraging heuristic and supervised functions from (Almasian et al., 2022), we extend this approach to neural models. We point out the limitation of the disjoint ranking and propose a quantity-centric fine-tuning paradigm for neural IR systems for the joint ranking of quantity and textual content.

3.1 Quantity Extraction

To facilitate both approaches, a prerequisite is a quantity extractor capable of identifying values (v), units (u), numerical conditions (c), and concepts (cn) associated with quantities. Concepts represent objects or events, and numerical values refer to. For instance, in the sentence "The iPhone 11 has 64GB of storage", the concept is "iPhone 11 storage". For this purpose, we use the Comprehensive Quantity Extractor (CQE) framework (Almasian et al., 2023). However, this module can be substituted with any alternative quantity extractor.

3.2 Disjoint Quantity Ranking

The disjoint model is based on the separation of quantity and term ranking. We assume that the

textual relevance of a sentence to query terms is 184 independent of the proximity of query and sentence 185 values under the query condition. Then, the relevance of a sentence can be the summation of (1)187 textual similarity, and (2) proximity of quantities in a sentence and a query, denoted in Eq 2. Note 189 that here, sim computes the similarity of search 190 terms to a sentence independent of sim_c , which 191 computes the quantity proximities given query con-192 dition c. τ and τ' signify that representations for 193 query and document are not necessary created from 194 the same model. If the query is not quantity-centric, 195 simply by removing the quantity score sim_c , the 196 models fall back to term scoring. 197

$$r(s_i|X) \sim sim(\tau(T_x), \tau(T_i)) + sim_c(\tau'(q_x), \tau'(Q_i))$$
(2)

In the following, we describe the computation of term (1), and quantity (2) scorings, where (1) $sim(\tau(T_x), \tau(T_i))$ and (2) $sim_c(\tau(q_x), \tau(Q_i))$. The general pipeline is depicted in App A.1.

3.2.1 Quantity Scoring

198

199

200

205

207

208

210

211

212

213

214

215

216

217

Using a quantity index containing explicit information about values and units in normalized form, we use heuristic functions to compute the proximity of query and sentence values based on different numerical conditions.

Index creation: Documents are split into sentences that are processed independently by CQE. CQE outputs standardized values, e.g., \$300 million is converted to \$300,000,000 and normalized units, e.g., kilometer per hour and km/h are mapped to the same unit. A quantity index with unit/value pairs is built from this output and resembles a lexical index. Here, each unique unit/value pair points to a list of sentences it occurs in.

Scoring functions: $sim_c(\tau(q_x), \tau(Q_i))$ is esti-218 mated by a scoring function qs that ranks the value 219 in a sentence based on the value in the query given 220 a numerical condition, where higher values indicate higher relevance. qs is dependent on the numerical condition, resulting in different scores for the same values under different numerical condi-224 tions. The quantity score only matters if the units 225 match, otherwise, the values are not comparable and refer to different aspects of an object, e.g., the horsepower of a car is different from the km/h it reaches. qs is formulated in Eq 3. The indicator function $\mathbb{1}_{u_i}(u_x)$ enforces the match between the units of the query and the sentence, and Φ_c is the condition-dependent scoring function. To obtain a value between 0 and 1, the score is normalized 233

by the number of quantities $|Q_i|$ in s_i . For brevity, from now on we refer to qs(s, c, X) as simply qs.

$$qs(s_i, c, X) := \frac{1}{|Q_i|} \sum_{i=1}^{|Q_i|} \mathbb{1}_{u_i}(u_x) \Phi_c(v_x, v_i) \quad (3)$$

 Φ_c consists of three heuristic functions, one for each numerical condition (*equal, less than, greater than*), adapted from (Almasian et al., 2022). The study in (Almasian et al., 2022) explores various Φ functions and their implications for sorting of results (Refer to App A.2). Simply by changing the Φ s, results can be rearranged, independent of training data that might introduce bias for a specific sorting preference. Nonetheless, for the evaluation of our model against other baselines, we focus only on the most intuitive variant, which ranks quantities with values closer to the query value in descending order. The Φ s are defined in Eq. 4. v_x is the query value and v_i is the sentence value.

$$\begin{split} \Phi_{=}(v_{x}, v_{i}) &=: exp(-|v_{x} - v_{i}|) \\ \Phi_{>}(v_{x}, v_{i}) &=: \begin{cases} v_{x}/v_{i} & v_{x} > v_{i} \\ 0 & else \end{cases} \tag{4} \\ \Phi_{<}(v_{x}, v_{i}) &=: \begin{cases} v_{i}/v_{x} & v_{x} < v_{i} \\ 0 & else \end{cases} \end{split}$$

 $\Phi_{=}$ assesses the proximity of v_x to v_i by employing the exponential decay of their difference. The resulting score ranges between 0 and 1, with larger absolute differences yielding lower scores.

The scoring functions $\Phi_{<}$ and $\Phi_{>}$ determine numerical proximity based on the ratio of the query value v_x to the sentence value v_i , resulting in a score between 0 and 1. This ratio, independent of magnitude, yields higher scores for closer values.

3.2.2 Term Scoring

Term scoring, $sim(\tau(T_x), \tau(T_i))$, can come from any lexical or semantic ranker, requiring only normalized scores. Yet, IR systems typically do not normalize their scores, as it has no influence on the final ranking. Here, we discuss ways to normalize scores of lexical and semantic systems and combine them with qs. For a lexical model we use BM25 (Robertson and Zaragoza, 2009) and for dense and sparse neural rankers, ColBERT (Khattab and Zaharia, 2020) and SPLADE (Formal et al., 2021) are employed.

Lexical model: Following (Almasian et al., 2022), we combine qs with the BM25 score. The combined score, represented in Eq 5, is constrained

253

254

255

256

257

258

259

260

261

262

263

264

265

266

267

269

270

271

272

273

274

275

234

235

236

237

238

239

240

241

242

243

244

245

246

247

248

249

to sentences containing the search terms, as indi-276 cated by $\mathbb{1}_{T_r}(s_i)$. The parameter α controls the influence of the quantity scoring, falling back to 278 pure term-based scoring when α is set to zero. The 279 BM25 (s_i, T_x) score is normalized per query by dividing each sentence's score by the maximum BM25 score retrieved for the specified search terms $\max_X = \max_{s \in S}(BM25(s, T_x)).$

277

281

289

290

294

297

301

307

311

314

$$QBM25(s_i, c, X) := \frac{BM25(s_i, T_x)}{\max_X} + \alpha \mathbb{1}_{T_x}(s_i)qs$$
(5)

Neural dense model: Representing a dense neural model, ColBERT is selected for term scoring. This choice is due to the same model being used for joint quantity ranking, where token-level interactions are crucial. Contextualized term score is computed with the similarity computation between token embeddings of query and sentence, as in Eq 6. ColBERT utilizes two BERT (Devlin et al., 2019) encoders for query and document (sentence), where each encoder outputs a list of token embeddings.

$$ColBERT(T_x, s_i) = \sum_{k \in |BERT(T_x)|} max_{j \in |BERT(s_i)|} BERT(T_x)_k \cdot BERT(s_i)_j$$
(6)

The term score comes from the MaxSim operation between the query and sentence embeddings. 298 MaxSim calculates an unbounded score for the maximum cosine similarity among the token em-299 beddings. To normalize this score, we require the maximum score. However, calculating the maximum score for the entire collection is impractical. For ranking, ColBERT leverages the pruningfriendly nature of the MaxSim in an approximate 304 nearest neighbor search (Johnson et al., 2019) to return top-k most relevant candidate sentences S_k . We compute the maximum score based on these candidate sentences $max_X = max_{s \in S_k}$ (ColBERT) to normalize the score between 0 and 1. qs is then exclusively applied to the top-k candidates, serving 310 as a second-stage re-ranker for numerical proximity. The final score is defined in Eq 7, α controls 312 the impact of quantity scoring. 313

$$\text{QColBERT}(s_i, c, X) := \frac{\text{ColBERT}(T_x, s_i)}{\max_X} \alpha \cdot qs$$
(7)

315 Note that the qs only affects the top-k sentences. 316 We also present a neural sparse model, where qs is integrated into the entire ranking. 317

Neural sparse model: The SPLADE model extends the document and query terms and uses an 319

inverted index for sparse dot products, allowing for end-to-end integration with the quantity scoring. Instead of term frequencies inside the index, term importance weights are computed by SPLADE. For each sentence and query, the BERT embeddings are passed through a ReLU non-linearity and log function to produce a sparse vector over the entire vocabulary, where the values inside this vector are the term importance. Then the relevance of the query to a sentence is based on the sparse dot product of this vector, as denoted in Eq 8.

320

321

322

323

324

325

326

327

329

330

331

333

334

335

336

337

339

340

341

342

343

344

345

346

347

348

349

351

352

353

354

355

356

357

360

361

362

363

364

365

$$SPLADE(s_i, T_x) := \log(1 + ReLU(BERT(s_i))) \cdot \log(1 + ReLU(BERT(T_x)))$$
(8)

We normalize the SPLADE score by the maximum score for a given query, \max_X = $max_{s_i \in S}$ (SPLADE (s_i, T_x)), as defined in Eq 9. For higher precision, quantity score is only added to sentences where there is a match between the expanded query terms and documents, denoted by the indicator function 1.

$$QSPLADE(s_i, c, X) := \frac{SPLADE(s_i, T_x)}{\max_X} + \alpha \mathbb{1}(s_i)qs$$
(9)

Joint Textual and Quantity Ranking 3.3

The independence assumption between the relevance of quantities and terms can be problematic. Consider the query "iPhone XR below €200". In a disjoint ranking, the following sentences can receive an inappropriately high score.

1) The price of an iPhone XR reached \in 236.50, whereas Samsung A14 is €132.00. This sentence has multiple quantities and the numerical condition is satisfied for a value unrelated to "iPhone XR".

2) Older iPhones, including iPhone XR have dropped in price with iPhone 8 to $\in 152.94$. Here, "iPhone XR" has no associated quantity.

These cases are due to a lack of correct association between concept and quantity. We refer to this as quantity-concept mismatch. To address this, we need to rank sentences based on quantities in context. Transformer-based models inherently capture token inter-dependencies across the entire context. However, current benchmarks lack quantity-centric data. Therefore, it remains unclear whether the deficiency in quantity understanding is due to the absence of task-specific training data or if the current architectures hinder numerical comparisons. To investigate this, we propose a data generation approach to address following problems.

460

461

410

First is the inability to perform value comparisons given numerical conditions. E.g., in the example above, the models ignore the *less than*, condition and focus on the semantic similarity of query text and sentence. Second, the semantic similarity of units is not well-defined. In the example above, results with "dollar" and other currencies receive high scores due to the context similarity of the units. Refer to App B.6 for a detailed discussion.

Our data generation paradigm is designed to enhance *value comparisons* and understanding of *unit surface forms*, by generating contrastive positive and negative sentences through data augmentation. Data augmentation, widely used in computer vision, has also found applications in NLP tasks (Sennrich et al., 2016). The GENBERT model (Geva et al., 2020) is a relevant example, which employs templates for generating pre-training data, to enhance numerical reasoning in question-answering systems without specialized architectures.

Similar to GENBERT, we fine-tune neural IR models used in the disjoint setting, ColBERT, and SPLADE, on synthetic data for quantity-centric IR¹. The data generation pipeline has three stages described in the following: *quantity extraction*, *query generation*, and *sample generation*.

3.3.1 Quantity Extraction

387

400

401

402

403

404

405

408

409

The documents are split into sentences and fed to CQE to extract quantities and concepts. The corpus is then transformed into an index-like structure based on concepts and units. We refer to this structure as *concept/unit index*. The keys of the index are concept/unit pairs that point to a list of values associated with the pair and a list of respective sentences they occur in. The list of values can be viewed as the distribution of values for a concept under a specific unit. An example entry is shown in App B.1. We utilize this index structures in the subsequent steps for query and sample generation.

3.3.2 Query Generation

406 For each concept/unit pair, three queries, one for 407 each condition, are created with the template

query = {concept} {numerical_condition}
{unit_before}{value}{unit_after}.

The variables enclosed in the brackets are populated during query generation. These steps are

given in an algorithm in the App B.2. In the following, we describe how each placeholder is filled.

Unit: A surface form of the query unit is chosen randomly from a dictionary of unit surfaces provided by CQE, e.g., " \in " is a surface of the unit "euro". unit_before and unit_after account for symbols appearing before, e.g., " \in " and abbreviations after a value, e.g., "EUR", respectively.

Value: For sample generation, sentences containing values meeting the query condition are crucial. Therefore, selecting query values with enough supporting sentences is vital. We propose the following strategy, based on the value distribution of *concept/unit index*:

Equal query: Query values are chosen from the most frequent values in the index (peak of value distributions), ensuring availability of maximum supporting sentences for a given concept and unit. *Less and greater than queries:* For these bounds, optimal candidates are close to the average of the value distribution, such that when the numerical condition is applied more sentence fall within limits. Infrequent values (tail of the distribution) may have inadequate supporting sentences for the sample generation step. Refer to the App B.3 for examples on value selection.

To avoid systemic bias by focusing on the most frequent values, we generate a second set of queries for each unit and concept pair by picking the query values at random.

To account for variability in representation, surface forms of large values that have multiple written forms are randomly replaced with their written form. This takes the shape of a composite of numbers and postfixes, such as "10 million," or includes commas for digit separation, e.g., "10,000,000".

Numerical Condition: This is a phrase in natural language indicating a bound on a quantity, e.g., "above" for *greater than* condition. For this purpose, a surface-form dictionary is created, and the respective placeholder is filled with values randomly chosen from the dictionary (see App B.4).

Concept: CQE identifies multi-word spans in a sentence as concepts. Utilizing them directly for query generation overlooks the nuances of semantic queries. For example, in the sentence "Disney+ charges \$6.99 a month.", "Disney+" is the extracted concept. "Disney+" is a streaming platform, including other media services. Such a sentence is relevant for a lexical query with exact matches, e.g., "Disney+ price under \$7.99 a month", or for a semantic query, e.g., "streaming platform price over

¹Given that we are perturbing values and units in a sentence, one might alternatively call this *data perturbation*.

5 dollar/month". Relying exclusively on keywords 462 in sentences poses a risk of biasing the neural mod-463 els toward lexical search and away from semantic 464 search. To avoid such a case, we add concept ex-465 pansion, where a large language model, namely 466 GPT-3 (Chen et al., 2023), is used to generate syn-467 onyms or synsets for a given concept (see App B.5). 468 These expansions are used to generate semantic 469 queries. E.g., "Disney+" becomes "Streaming plat-470 form". For each expanded concept new values and 471 unit surface forms are sampled to generate seman-472 tic queries for each numerical condition. 473

3.3.3 Sample Generation

474

484

485

486

487

494

495

496

497

498

499

500

503

504

505

506

507

508

510

511

The input of this stage are the generated queries and 475 the concept/unit index. The sample generation step 476 477 creates positive and negative training samples for each quantity-centric query. This includes positive 478 and negative samples obtained directly from the 479 dataset, as well as additional augmented samples. 480 An overview of the sample generation pipeline and 481 482 an algorithmic view is presented in App B.7. In the following, we describe each step in detail. 483

Look-up: Given a query containing a (*concept*, *unit*, *condition*, *value*), we conduct a lookup in the *concept/unit index* to retrieve the sentences and the distribution of values.

488 **Positive and Negative Sentences List:** The ob-489 tained sentences are divided into positive s_+ and 490 negative s_- lists, based on the numerical condition. 491 s_+ contains sentences, where the values in them 492 satisfying the condition and s_- contains sentence 493 violating the condition.

Original sampling: With sample size n, sentences are randomly selected from s_+ as positive samples (s_{o+}) and from s_- as negative samples (s_{o-}) . Refer to App B.9 for information on the sample size.

Unit permutation sampling: This method generates positive and negative samples to cover diverse unit surface forms using CQE's unit dictionary. Positive samples contain various surface forms of the unit in the query, while negative samples include surface forms of units in the same family as the query unit, creating negatives.

- A positive sample, s_{u+} , is formulated by substituting the unit in a positive sentences, s_+ , with other surface forms of unit in query u_i .
- A negative sample, s_u is created by replacing the unit in a positive sentences, s₊, with a surface form of a unit different from query unit ,u_i, but belonging to the same family. The

unit families are grouping based on the property they measure. For example, "pace", "meter", and "foot" all belong to the family of "length". Sampling the surface form from the same family ensures a fine distinction between unit types, even in similar contexts.

512

513

514

515

516

517

518

519

520

521

522

523

524

525

526

527

528

529

530

531

532

533

534

535

536

537

538

539

540

541

542

543

544

545

546

547

548

549

550

551

552

553

554

555

556

557

558

Value permutation sampling: This permutation emphasizes the importance of the value comparison and numerical conditions, highlighting that sentence relevance depends on whether the sentence value satisfies the query condition or not.

- A positive sample, s_{v+}, is formulated by permuting the values in a negative sentence s₋, maintaining the correct concept and unit but adjusting the value to satisfy the quantity condition.
- A negative sample, s_{v-}, is generated by permuting the values in a positive sentence s₊, where concept, unit, and value are all correct, to invalidates the quantity condition.

The replacement values are sampled from the values in the *concept/unit index*, mirroring the underlying distribution of the relevant quantity, as to the reason for this choice, refer to the App B.8.

Aggregate: The final set of positive and negative samples for each query is the union of all samples generated from the original sampling, value and unit permutation, $s_{f+} = s_{o+} \cup s_{u+} \cup s_{v+}$ and $s_{f-} = s_{o-} \cup s_{u-} \cup s_{v-}$.

The models reported in the evaluation use a combination of original sampling with unit permutation and concept expansion on the query. Value permutation did not show stable performance gains, which we attribute to the difficulty of numerical representations in dense models. For more discussion on this matter and ablation study of augmentation methods refer to App C.5.

4 Evaluation

Given the absence of task-specific models, we assess our quantity-aware models against general domain lexical and neural models.

Lexical models include a BM25 and a BM25_{*filter*} variant. BM25_{*filter*} has a separate numerical index to eliminate the results of BM25 where the query condition is not met. This method resembles the numerical indices from databases, focusing on filtering rather than ranking.

Table 1: Query types in FinQuant and MedQuant.

| | FinQuant | MedQuant |
|-----------------------|----------|----------|
| Total queries | 420 | 210 |
| Sentence in corpus | 306,291 | 153,252 |
| Per condition | 140 | 70 |
| Keyword-based queries | 300 | 120 |
| Semantic queries | 120 | 90 |

Neural models include the trained checkpoints of SPLADE and ColBERT as well as Cohere_{v3}². Cohere_{v3} embeddings are included to show that even industry-level models trained on extensive data still lack quantity-centric understanding.

4.1 Datasets

559

560

561

562

564

565

566

570

571

574

575

576

581

583

586

589

591

594

595

We introduce two English resources called Fin-Quant and MedQuant. To the best of our knowledge, these are the first quantity-centric benchmarks for retrieval. Test queries were manually formulated using the concept/unit index, covering both lexical and semantic queries. Statistics for various query types are presented in Table 1. There is an equal number of queries for each condition, and semantic queries constitute a smaller portion due to annotation challenges. For details on the dataset creation, refer to App C.1. The data is annotated by the two authors of the paper, with inter-annotator agreement computed on a subset of 20 samples per dataset. The Cohen's Kappa coefficient (Cohen, 1960) is 0.83 and 0.88 for FinQuant and MedQuant, respectively. FinOuant corpus contains over 300k sentences from 473,375 news articles. MedQuant is smaller, containing over 150k sentences from 375,580 medical documents of the TREC Medical Records (Voorhees, 2013). Since the concept/unit index is used for dataset creation, CQE's performance directly affects the data quality. While CQE is adept at handling financial data, extractions on clinical data were noisy, impacting performance comparisons later on. However, we find it important to report results on both datasets, making the reader aware of the lesser quality of MedQuant.

4.2 Ranking Performance

Table 2 shows the ranking performance of quantityaware models, in terms of P@10, MRR@10, NDCG@10, R@100, and latency in milliseconds. The three models with a "Q" prefix indicate the disjoint and unsupervised rankers. Neural models with a f_t postfix are joint models fine-tuned on synthetic data. Permutation re-sampling is used to test for significant improvements (Riezler and Maxwell, 2005). Results denoted with \dagger mark highly significant improvements over the baseline models, without quantity awareness with a *p*-value < 0.01. All results are from single runs. For Implementation details refer to App C.3. 599

600

601

602

603

604

605

606

607

608

609

610

611

612

613

614

615

616

617

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

646

647

648

649

Contrary to our initial hypothesis, disjoint rankers consistently outperform joint models across all metrics, with improvements exceeding 10 points in P@10 and over 30 points in MRR and NDCG over the base models (without the "Q" prefix), without requiring additional fine-tuning. The only drawback of disjoint models is a minimal increase in latency, especially for QBM25 and QSPLADE, where the quantity score is added to the entire ranking. This overhead diminishes for the topperforming model, QColBERT, where the quantity score serves as a re-ranker on the top-k candidates. ColBERT shows high recall on both datasets, suggesting that relevant results are within the top-k but not necessarily at the very top. Hence, the reranking with the quantity score proves beneficial. Joint models show a comparable performance boost, with metrics falling below those of the disjoint ranker but still improving from the base models. This validates our hypothesis that the absence of task-specific data has amplified the challenge of quantity understanding for retrieval systems. Here, once again the $ColBERT_{ft}$ variant shows superior performance. We attribute the better performance of the ColBERT-based model to the fine-grained token-level interactions that allow the model to learn better associations between tokens. In quantity ranking, token interactions play a more significant role compared to the query and document expansions conducted by SPLADE. This also showcases that the architecture and how the inter-token interactions are modeled matter for quantity understanding. Nonetheless, even after fine-tuning, understanding numerical conditions remains a challenge. We investigate how much the fine-tuned models rely on quantities for ranking in App C.4.

4.3 Cross-dataset Generalization

Two lower bottom rows of Table 2 list the performance drop of joint rankers on out-of-domain data, compared to models fine-tuned on generated data from the same domain. Each model is fine-tuned on data from the other dataset and shows minimal performance drop, suggesting that the models learn patterns for quantity-centric queries without mem-

²https://cohere.com/embeddings DLA: 10.02.24

| | Model | Model latency | | FinQuant | | | MedQuant | | | |
|---------------|------------------------|----------------|-------------------|-------------------|-------------------|------------------|-------------------|-------------------|-------------------|-------------------|
| | Widder | (ms) | P@10 | MRR@10 | NDCG@10 | R@100 | P@10 | MRR@10 | NDCG@10 | R@100 |
| | BM25 | 9 | 0.06 | 0.14 | 0.09 | 0.47 | 0.04 | 0.11 | 0.07 | 0.37 |
| baselines | BM25 filter | 9 | 0.14 | 0.32 | 0.25 | 0.60 | 0.08 | 0.19 | 0.15 | 0.48 |
| | $Cohere_{v3}$ | - | 0.14 | 0.22 | 0.19 | 0.27 | 0.10 | 0.17 | 0.15 | 0.25 |
| | SPLADE | 26 | 0.10 | 0.24 | 0.19 | 0.53 | 0.11 | 0.25 | 0.20 | 0.58 |
| | ColBERT | 36 | 0.15 | 0.35 | 0.27 | 0.70 | 0.12 | 0.31 | 0.24 | 0.63 |
| | QBM25 | 311 | 0.21 | 0.53 | 0.41 | 0.55 | 0.18 | 0.47 | 0.37 | 0.51 |
| joint | QSPLADE | 319 | 0.29^\dagger | 0.67^\dagger | 0.53^\dagger | 0.83^\dagger | 0.19^{\dagger} | 0.52^\dagger | 0.38^{\dagger} | 0.70^{\dagger} |
| | QColBERT | 42 | 0.30^\dagger | 0.69^\dagger | 0.56^\dagger | 0.87^\dagger | 0.18^{\dagger} | 0.51^\dagger | 0.37^\dagger | 0.73^\dagger |
| dista ind | $SPLADE_{ft}$ | 26 | 0.21 [†] | 0.51 [†] | 0.41 [†] | 0.74^{\dagger} | 0.14 [†] | 0.37 [†] | 0.29 [†] | 0.63 [†] |
| disjoint | $ColBERT_{ft}$ | 36^{\dagger} | 0.23^{\dagger} | 0.55^{+} | 0.44^{\dagger} | 0.77^{\dagger} | 0.18^{\dagger} | 0.44^{\dagger} | 0.36 [†] | 0.72^\dagger |
| anaaa dataaat | SPLADE _{out} | 26 | -0.03 | -0.06 | -0.07 | -0.04 | -0.02 | -0.01 | -0.04 | -0.05 |
| cross-dataset | ColBERT _{out} | 36 | -0.03 | -0.07 | -0.06 | -0.03 | -0.02 | -0.01 | -0.03 | -0.02 |

Table 2: P@10, MRR@10, NDCG@10 and R@100 for on FinQuant and MedQuant. Top-2 values in each column are highlighted in bold.

orizing common queries.

4.4 Lexical vs Semantic Queries

Fig 1a shows NDCG@10 of all models on lexical and semantic subsets of the FinQuant. The *seen* and *unseen* are lexical queries and *expansion* and *w/o surface form* represent semantic queries. For the details of their distinction refer to App C.2. Interestingly, the disjoint ranking using dense models captures both semantic similarity and quantity understanding. QBM25 performs equally well as dense models in lexical queries but significantly worse on semantic ones. Joint rankers outperform base models in both lexical and semantic queries but lag behind disjoint models.

Fig 1b depicts NDCG@10 of all models on different numerical conditions. *Equal* queries are in general easier for the models as the notion of relevance in this case aligns with textual ranking. The performance drops almost 20 points for the boundbased conditions. This drop is consistent across all models, implying that the bound-based conditions are harder for models to rank.

5 Conclusion and Ongoing Work

This work introduces two methods to integrate 673 quantity understanding into existing retrieval sys-674 tems. The disjoint approach is an unsupervised and heuristic method, while the joint approach involves 676 fine-tuning on quantity-centric synthetic data to 677 enhance quantity understanding. The disjoint scoring can be combined with any lexical or semantic matchers without the change in their architecture or need for fine-tuning, showing consistently good 681 performance regardless of data distribution. Moreover, the notation of quantity proximity is easily altered by changing the quantity scoring function,



Figure 1: Performance on different subsets of FinQuant.

685

686

687

688

689

690

691

692

693

694

695

696

697

698

699

leading to great flexibility in terms of different sorting of results. However, the quantity index introduces an overhead in query latency, is sensitive to errors from quantity extraction and the independence assumption leads to possible conceptquantity mismatch in the results. Conversely, the joint models are better at finding concept and quantity associations but their overall performance is lower. Yet, if one does not want to be dependent on an external index and quantity extractor, the fine-tuning on synthetic data can enhance quantity awareness to some extend. We also introduce two benchmark datasets and evaluated our methods against multiple baselines. In future, we explore the impact of numerical embeddings in retrieval.

652 653

6 Limitations

In this section, we highlight the limitations of the
proposed evaluation resources and the models introduced in this paper.

Evaluation resources: One immediate consideration about the datasets is the relatively limited number of test queries compared to larger-scale datasets such as MSMARCO (Nguyen et al., 2016). 707 This is mainly due to limited human resources and budget in an academic setting. Nonetheless, we argue that this number of the query is already enough 710 to showcase certain quantity-centric capabilities. 711 Another shortcoming of the data is the absence of 712 queries for ranges, e.g., "iPhone with price between 713 500 and 800 dollars", and negations, "iPhones not 714 equal to 500 dollars". 715

Quantity-aware models: When considering neu-716 ral models, one limitation is their reliance on hard-717 ware capabilities, particularly the need for GPUs 718 to ensure efficient training, indexing, and inference. 719 The query latency values reported in this paper would suffer greatly if the computations were don 721 on CPU. Moreover, both the synthetic data gener-722 ation paradigm and the disjoint model rely on a 723 quantity extractor. In the case of the disjoint model, the quality of the quantity index directly relies on the quality of value and unit extraction. If a value 726 and unit is not detected by the extractor it will not 727 be considered by the scoring function. In the joint model, for data generation, the quantity extractor should also possess the ability to detect concepts in text, introducing the potential for additional error 731 propagation through the system. In this work, we do not discuss models that deal with ranges and negations. Adding such variations to the disjoint 734 models requires only a change in the numerical 735 scoring function but it is more difficult for the joint setting where proper training data is required. For the bound-based conditions of *less than* and *greater* 738 than, we considered open bounds. Depending on 739 the user intent closed bounds might be more appro-740 priate, however, similar to the optimal sorting of 741 results, this issue does not have a single solution.

References

743

744

745

746

747

748

- Rakesh Agrawal and Ramakrishnan Srikant. 2003. Searching with Numbers. *IEEE Trans. Knowl. Data Eng.*, 15(4):855–870.
 - Satya Almasian, Milena Bruseva, and Michael Gertz. 2022. QFinder: A Framework for Quantity-centric

Ranking. In SIGIR '22: The 45th International ACM SIGIR Conference on Research and Development in Information Retrieval, pages 3272–3277. ACM.

- Satya Almasian, Vivian Kazakova, Philip Göldner, and Michael Gertz. 2023. CQE: A Comprehensive Quantity Extractor. In *Proceedings of the 2023 Conference on Empirical Methods in Natural Language Processing, EMNLP*, pages 12845–12859. ACL.
- Zekai Chen, Mariann Micsinai Balan, and Kevin Brown. 2023. Language Models are Few-shot Learners for Prognostic Prediction. *CoRR*, abs/2302.12692.
- Jacob Cohen. 1960. A Coefficient of Agreement for Nominal Scales. *Educational and Psychological Measurement*, 20(1):37–46.
- Jacob Devlin, Ming-Wei Chang, Kenton Lee, and Kristina Toutanova. 2019. BERT: Pre-training of Deep Bidirectional Transformers for Language Understanding. In Proceedings of the 2019 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, pages 4171–4186. ACL.
- Marcus Fontoura, Ronny Lempel, Runping Qi, and Jason Y. Zien. 2007. Inverted Index Support for Numeric Search. *Internet Math.*, 3(2):153–185.
- Thibault Formal, Benjamin Piwowarski, and Stéphane Clinchant. 2021. SPLADE: Sparse Lexical and Expansion Model for First Stage Ranking. In *The 44th International ACM SIGIR Conference on Research and Development in Information Retrieval*, pages 2288–2292. ACM.
- Mor Geva, Ankit Gupta, and Jonathan Berant. 2020. Injecting Numerical Reasoning Skills into Language Models. In *Proceedings of the 58th Annual Meeting of the Association for Computational Linguistics, ACL*, pages 946–958. ACL.
- Vinh Thinh Ho, Yusra Ibrahim, Koninika Pal, Klaus Berberich, and Gerhard Weikum. 2019. Qsearch: Answering Quantity Queries from Text. In *The Semantic Web - ISWC - 18th International Semantic Web Conference, Proceedings*, volume 11778 of *Lecture Notes in Computer Science*, pages 237–257. Springer.
- Vinh Thinh Ho, Koninika Pal, Niko Kleer, Klaus Berberich, and Gerhard Weikum. 2020. Entities with Quantities: Extraction, Search, and Ranking. In WSDM '20: The Thirteenth ACM International Conference on Web Search and Data Mining, pages 833–836. ACM.
- Yusra Ibrahim, Mirek Riedewald, Gerhard Weikum, and Demetrios Zeinalipour-Yazti. 2019. Bridging Quantities in Tables and Text. In *35th IEEE International Conference on Data Engineering, ICDE*, pages 1010– 1021. IEEE.
- Chengyue Jiang, Zhonglin Nian, Kaihao Guo, Shanbo Chu, Yinggong Zhao, Libin Shen, and Kewei Tu. 2020. Learning Numeral Embedding. In *Findings*

749

750

763

764

765

766

767

768

769

770

771

773

774

775

776

777

778

779

783

784

785

786

787

788

789

790

791

792

793

794

795

796

797

798

799

800

801

802

854 855

of the Association for Computational Linguistics: EMNLP 2020, Online Event, 16-20 November 2020, volume EMNLP 2020 of Findings of ACL, pages 2586-2599. ACL.

- Zhihua Jin, Xin Jiang, Xingbo Wang, Qun Liu, Yong Wang, Xiaozhe Ren, and Huamin Qu. 2021. NumGPT: Improving Numeracy Ability of Generative Pre-trained Models. CoRR, abs/2109.03137.
- Jeff Johnson, Matthijs Douze, and Hervé Jégou. 2019. Billion-scale similarity search with GPUs. IEEE Transactions on Big Data, 7(3):535–547.
- Omar Khattab and Matei Zaharia. 2020. ColBERT: Efficient and Effective Passage Search via Contextualized Late Interaction over BERT. In Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, pages 39-48. ACM.
- Tongliang Li, Lei Fang, Jian-Guang Lou, Zhoujun Li, and Dongmei Zhang. 2021. AnaSearch: Extract, Retrieve and Visualize Structured Results from Unstructured Text for Analytical Queries. In WSDM '21, The Fourteenth ACM International Conference on Web Search and Data Mining, pages 906–909. ACM.
- Arun S. Maiya, Dale Visser, and Andrew Wan. 2015. Mining Measured Information from Text. In Proceedings of the 38th International SIGIR Conference on Research and Development in Information Retrieval, pages 899-902. ACM.

Tri Nguyen, Mir Rosenberg, Xia Song, Jianfeng Gao, Saurabh Tiwary, Rangan Majumder, and Li Deng. 2016. MS MARCO: A Human Generated MAchine Reading COmprehension Dataset. In Proceedings of the Workshop on Cognitive Computation: Integrating neural and symbolic approaches 2016 co-located with the 30th Annual Conference on Neural Information Processing Systems (NIPS), volume 1773 of CEUR Workshop Proceedings. CEUR-WS.org.

- Rodrigo Frassetto Nogueira, Zhiying Jiang, and Jimmy Lin. 2021. Investigating the Limitations of the Transformers with Simple Arithmetic Tasks. CoRR, abs/2102.13019.
- Oiu Ran, Yankai Lin, Peng Li, Jie Zhou, and Zhiyuan Liu. 2019. NumNet: Machine Reading Comprehension with Numerical Reasoning. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing EMNLP-IJCNLP, pages 2474-2484. ACL.
- Stefan Riezler and John T. Maxwell. 2005. On some pitfalls in automatic evaluation and significance testing for MT. In Proceedings of the ACL Workshop on Intrinsic and Extrinsic Evaluation Measures for Machine Translation and/or Summarization, pages 57-64, Ann Arbor, Michigan. ACL.
- Stephen Robertson and Hugo Zaragoza. 2009. The Probabilistic Relevance Framework: BM25 and Beyond. Now Publishers Inc.

Maciej Rybinski, Stephen Wan, Sarvnaz Karimi, Cécile Paris, Brian Jin, Neil I. Huth, Peter J. Thorburn, and Dean P. Holzworth. 2023. SciHarvester: Searching Scientific Documents for Numerical Values. In Proceedings of the 46th International ACM SIGIR Conference on Research and Development in Information Retrieval, SIGIR, pages 3135–3139. ACM.

861

862

864

865

866

867

868

869

870

871

872

873

874

875

876

877

878

879

880

881

882

883

884

885

886

888

889

890

891

892

893

894

895

896

897

899

900

901

902

903

904

905

906

907

908

909

910

911

912

913

914

915

- Rico Sennrich, Barry Haddow, and Alexandra Birch. 2016. Improving Neural Machine Translation Models with Monolingual Data. In Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics, ACL. ACL.
- Georgios P. Spithourakis and Sebastian Riedel. 2018. Numeracy for Language Models: Evaluating and Improving their Ability to Predict Numbers. In Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics, ACL 2018, Melbourne, Australia, July 15-20, 2018, Volume 1: Long Papers, pages 2104-2115. ACL.
- Daniel Spokoyny, Ivan Lee, Zhao Jin, and Taylor Berg-Kirkpatrick. 2022. Masked Measurement Prediction: Learning to Jointly Predict Quantities and Units from Textual Context. In Findings of the Association for Computational Linguistics: NAACL 2022, pages 17-29. ACL.
- Dhanasekar Sundararaman, Shijing Si, Vivek Subramanian, Guoyin Wang, Devamanyu Hazarika, and Lawrence Carin. 2020. Methods for Numeracy-Preserving Word Embeddings. In Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing, EMNLP, pages 4742-4753. ACL.
- Avijit Thawani, Jay Pujara, Filip Ilievski, and Pedro A. Szekely. 2021. Representing Numbers in NLP: a Survey and a Vision. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, NAACL-HLT, pages 644-656. ACL.
- Christophe Van Gysel and Maarten de Rijke. 2018. Pytrec_eval: An Extremely Fast Python Interface to trec_eval. In SIGIR. ACM.
- Ellen M. Voorhees. 2013. The TREC Medical Records Track. In ACM Conference on Bioinformatics, Computational Biology and Biomedical Informatics. ACM-BCB 2013, page 239. ACM.
- Eric Wallace, Yizhong Wang, Sujian Li, Sameer Singh, and Matt Gardner. 2019. Do NLP Models Know Numbers? Probing Numeracy in Embeddings. In Proceedings of the 2019 Conference on Empirical Methods in Natural Language Processing and the 9th International Joint Conference on Natural Language Processing, EMNLP-IJCNLP, pages 5306-5314. ACL.
- Gerhard Weikum. 2020. Entities with Quantities. IEEE Data Eng. Bull., 43(1):4–8.



Figure 2: General pipeline of the disjoint quantityranking approach, where a separate quantity index is responsible for computing quantity proximity and a termbased lexical or semantic index is used to compute the similarity of the search terms to sentences.

A Disjoint Quantity Ranking

In this section, we provide additional material related to the disjoint quantity ranking model.

A.1 Disjoint Quantity Ranking Pipeline

The general pipeline for the disjoint quantityranking model is shown in Figure 2. The query is processed into quantity, search terms, and conditions, using CQE or similar packages. The document corpus is indexed separately for terms and quantities, whereas the term-based index can be a traditional lexical index or a vector database. This term-based index retrieves semantically or lexically similar sentences. From the retrieved sentences, the quantity index identifies values that share the same unit as the query, computing proximity based on the provided condition. The final ranking combines scores from term-based and quantity ranking.

A.2 Optimal Sorting

Although all the sentences that satisfy a numerical condition and have the correct concept and unit are potentially relevant, the order in which the result items are presented to the user can either aid or hinder the user in finding the desired result. In term-based ranking, the optimal order of results is evident. However, when it comes to quantities, relevance is more subjective and the optimal sorting is dependent on the user's information needs. For example, a user searching for "iPhone camera that has more than 8 inches" might look for a maximum value larger than 8 inches or a display only marginally larger, both of which are valid answers. Presenting results in ascending or descending order based on numerical distances allows the user to identify the desired result more efficiently. (Almasian et al., 2022) briefly addresses this issue

and explores potential alternatives for scoring functions to enable various sorting options. Disjoint approach is flexible concerning different sorting. By switching a scoring function, the results can be rearranged. Joint model are not as adaptive, and rearranging the results requires additional fine-tuning based on a new preferred sorting. 952

953

954

955

956

957

958

959

960

961

962

963

964

965

966

967

968

969

970

971

972

973

974

975

976

977

978

979

980

981

982

983

984

985

986

987

988

989

B Joint Quantity Ranking

In this section, we provide the additional material related to the joint quantity ranking model.

B.1 Concept-unit Index

An entry in the *concept/unit index* from the Fin-Quant dataset is shown below.

```
{("cannabis company","cent per share"):
{"values":[1.4, 17.0, 17.0, 22.0, 26.0,
35.0, 84.0],
"sentences":['The cannabis company says
the loss amounted to 0.9 of a cent per
share for the quarter ended May 31
compared with a loss of \$4 million or
1.4 cents per share a year earlier .',
'The cannabis company says its loss
amounted to 17 cents per diluted share
for the quarter ended Jan. 31 .',...]}}
```

Note that repetition of values for the same concept/unit pair is stored as duplicates, such that the frequency of values is kept for the distribution, e.g., value "17.0" is repeated twice as it occurs in two distinct sentences. The creation steps are depicted in Figure 3. The corpus is analyzed with CQE to extract values, units, and concepts from each sentence, where sentences sharing the same unit and concept are grouped into a list, along with values represented as a distribution.

B.2 Query Generation

The complete query generation pipeline is depicted in Figure 4. The *concept/unit index* is used to select values and units for numerical conditions. Additionally, a large language model is used to expand concepts for semantic queries. The template generation block combines all the outputs to formulate three queries for each unit and concept pair. To generate expanded concepts, a new query value is chosen from the value distribution, leading to the formulation of a new set of queries. Additionally, we offer the query generation pseudocode in Algorithm 1 to make the input and output of each step clear. In the algorithm, v refers to the value, u to the unit, c to the condition, and cn to the concept.

927

929

930

931

932

935

936

938

917

918

919

943

947

948



Figure 3: Overview of the quantity tagging step and creation of *concept/unit index* structure.



Figure 4: Overview of the query generation pipeline, using *concept/unit* index and a large language model for concept expansion.



Figure 5: An example of choosing query values for *equal* and bound-based conditions.

| Algorithm | 1 (| Juery | Generation |
|-----------|-----|-------|------------|
|-----------|-----|-------|------------|

| function GENERATE_QUERY (cn, u, c) |
|--|
| $v \leftarrow get_query_values(cn_unit_dict, c)$ |
| $u_b, u_a \leftarrow \text{get_unit_surfaceform}(u)$ |
| $c \leftarrow \text{get_condition_surfaceforms}(c)$ |
| $query \leftarrow conc + c + u_b + v + u_a$ |
| return query |
| end function |
| |
| an whit dist (soneant/whit index |

 $\begin{array}{l} {\rm cn_unit_dict} \leftarrow {\rm concept/unit~index} \\ {\rm cn_expand_dict} \leftarrow {\rm concept~expansions} \\ {\rm for~}(cn,u) \mbox{ in cn_unit_dict~do:} \\ {\rm for~}cn \mbox{ in } [cn, {\rm cn_expand_dict}[cn]] \mbox{ do:} \\ {\rm for~}c \mbox{ in } (equal, greater, less) \mbox{ do:} \\ {\rm GENEARATE_QUERY}(cn,u,c) \\ {\rm end~for} \\ {\rm end~for} \\ {\rm end~for} \end{array}$

B.3 Choosing the Right Query Value

Each entry in the concept/unit index points to the 991 sentences and list of values in those sentences. For 992 the data augmentation to work, we require a num-993 ber of positive and negative samples per query and 994 therefore, it is important to choose the value of the 995 query such that supporting sentences in the corpus 996 are present. A hypothetical example of value distri-997 bution is shown in Figure 5. For the *equal* query, 998 the challenge is to find enough positive samples, 999 since there is an abundance of not equal values in each distribution. In Figure 5, values with the high-1001 est frequency, denoted by red arrows pointing to 1002 peaks in the distribution, serve as optimal candi-1003 dates for the equal condition. In this manner, we 1004 make sure that there a enough positive samples for 1005 the data augmentation. Values close to the average (highlighted in a yellow box) are chosen for 1007 the Less than and greater than queries. For such 1008 queries, we avoid infrequent values towards the tail of the distribution, to avoid too few positive or 1010 negative samples. 1011

B.4 Dictionary of Numerical Conditions

A non-comprehensive dictionary of surface forms1013for numerical conditions is shown Table 3, con-
taining multiple surface forms for each condition.1014

Table 3: Numerical conditions used for query generation and their surface forms.

| Condition | Surface forms |
|--------------|---|
| Equal | exactly, exact, equals, equals to, for, with, of, at |
| greater than | greater than, more than, above, larger than, over, higher than, exceed, exceeding |
| Less than | smaller than, below, less than, fewer than, no more than, beneath |

1016

1017

1012

990

B.5 Concept Expansion

For concept expansion, we use the OpenAI API 31018and employ the text-davinci-003 model with1019few-shot learning. We set the temperature to 1 to1020encourage creative responses. Since the concepts1021come from two distinct domains of finance and1022

³https://openai.com/ DLA:11.02.2024

medicine, the few-shot examples vary accordingly. Below we specify the two prompts used for concept expansion, the result is stored in a concept expansion dictionary and utilized during query generation. The place-holder concept is replaced with a concept from the *concept/unit index*.

For finance domain:

Complete with the words super set or synonym, but do not reuse the exact same words, the word "Super Set" should not be in the response and response should have at least two words:

S&P 500 = stock market index Audi = car Oil prices = petroleum prices unemployment rate = unemployment percentage iPhone sales = phone sales Netflix shares = stock shares President Trump = President iPhone 11= iphone Hong Kong = citystake PEXA = Property Exchange Australia shares

 $\{concept\} =$

For medical domain:

Complete with the words super set or synonym, but do not reuse the exact same words, the word "Super Set" should not be in the response and response should have at least two words:

ophthalmic solution = eve medication Control group = treatment group irinotecan hydrochloride = chemotherapy drug monoclonal antibody = substitute antibodies MRI scans = Magnetic resonance imaging influenza H1N1 vaccine = flu vaccine HAI antibody response = Influenza-specific antibody response

{concept} =

B.6 When Semantic Search Back-fires

Semantic retrieval systems consider an entire context to find a fuzzy relevance to a query at hand. Often, this aligns well with the user's expectations. For instance, when searching for a "dark color evening dress", any dress that can be worn as an evening gown and has a dark color would be suitable. But as soon as the user becomes more



Figure 6: Overview of the sample generation using value and unit permutation.

1040

1041

1042

1043

1045

1046

1047

1048

1049

1050

1051

1052

1054

1055

1056

1057

1058

1060

1062

1063

1064

1065

1070

1075

specific like "blue evening dress", the embedding space could also bring a similar color like "teal" into the search result. Depending on the user's flexibility regarding the dress color, this behavior may or may not be desirable. Such hard constraints are challenging for neural models. Quantity-centric queries compose hard constraints on values and units where the fuzzy matching of context might do more harm than good. For instance, when searching for a "car with more than 320 hp", if the results contain a car with "360 brake horsepower" instead of horsepower the result is irrelevant. Both horsepower and brake horsepower are used in similar contexts but refer to different attributes. Horsepower measures the power generated by the engine, while brake horsepower measures how much of the power produced by the engine is sent to the wheels which makes the car accelerate. Another common problem is with currencies. Given that monetary values often appear in similar contexts, it becomes challenging for the neural models to differentiate between various currency units. The same applies to hard constraints on values, where based on a given numerical condition, values outside of that bound are considered irrelevant.

B.7 Sample Generation with Permutations

An outline of sample generation pipeline is shown 1066 in Figure 6. The input of this stage is the gener-1067 ated queries and the *concept/unit index*. For each 1068 quantity-centric query, a list of positive and negative samples are created by applying the numerical condition on the list of sentences from the index. 1071 The original positive and negative samples are then 1072 chosen at random from such a list. The same list 1073 is utilized as seed samples for data augmentation. 1074 Unit and value permutation are employed to generate augmented positive and hard negative samples. 1076 Hard negative are positive samples, where the unit or value is perturb to violet the query condition. 1078

1031

1023

1024

1025

1026

1027

1028 1029

1030

1035

The steps are presented in Algorithm 2. Each sampling mechanism is encapsulated within a distinct function and the final training samples are the union of all generation mechanisms. In the algorithm, v refers to value, *vals* to list of value of a given concept and unit, u to unit, c to condition, cn to concept and n to sample size.

| Algorit | hm 2 Sample Generation |
|---|--|
| functi end fu | on ORIGINAL_SAMPLING (s_+, s, n) return sample (s_+, n) , sample (s, n) unction |
| functi end fu | on UNIT_PERMUTATION (s_+, n, u) $s_{u+} \leftarrow \text{replace_same_unit_surface}(s_+, u)$ $s_{u-} \leftarrow \text{replace_other_unit_surface}(s_+, u)$ return sample (s_{u+}, n) , sample (s_{u-}, n) inction |
| functi end fu | on V_PERMUATION $(s_+, s, n, vals, c)$ $s_{u+} \leftarrow$ replace_with_positive_value (s, v) $s_{u-} \leftarrow$ replace_with_negative_value (s_+, v, c) return sample (s_{v+}, n) , sample (s_{v-}, n) inction |
| conc_u querie $n \leftarrow n$ for (cr | unit_dict \leftarrow concept/unit index s \leftarrow list of queries umber of samples n, u, c, v in queries do : |
| | $s, vals \leftarrow \text{conc_unit_dict}[(cn, u)]$ $s_{+}, s_{-} \leftarrow \text{filter_based_on_condition}(s, c, v)$ $s_{o+}, s_{o-} \leftarrow \text{ORIGINAL_SAMPLING}(s_{+}, s_{-}, n)$ $s_{u+}, s_{u-} \leftarrow \text{UNIT_PERMUTATION}(s_{+}, n, u)$ $s_{v+}, s_{v-} \leftarrow \text{V_PERMUTATION}(s_{+}, s_{-}, n, vals, v, c)$ |

end for

1079

1080

1081

1082

1083

1084

1085

1086

1087

1088

1089

1090

1092

1093

1094

1095

1096

1097

1098

1099

1100

1101

1102

1103

1104

B.8 Sampling within Distribution

It is crucial that the permuted values obey the original value distribution of the corpus. The properties of concepts are often limited to a specific range, e.g., the value "10000" is unreasonable for percentage rate of unemployment. Moreover, certain values are on a discrete scale with limited options, e.g., "RAM of a laptop" is limited to distinct values such as 4,8, and 16. Assigning a random number outside this range, like 10, would be unrealistic. Therefore, for the synthetic data to obey the rule of the real-world dataset and reflect the distribution of different properties, the permuted value are chosen from the values observed in the corpus.

B.9 Down-sampling

If the number of available sentences in the positive and negative lists is smaller than the sample size, a *downsampling* procedure is implemented. When $|s_+| < n$ or $|s_-| < n$, we reduce the sample size to the smallest number of available samples. 1105

1106

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124 1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

1149

1150

1151

1152

1153

C Evaluation

In this section, we present additional evaluations 1107 and implementation details. To reproduce the results, access model checkpoints and datasets, we encourage readers to refer to our repository. 1110

C.1 FinQuant and MedQuant Datasets

In this section, we give an overview of the creation the FinQuant and MedQuant evaluation benchmarks. FinQuant is created from a set of news articles in categories of economics, science, sports, and technology, collected between 2018 and 2022. MedQuant contains TREC Medical Records (Voorhees, 2013) on clinical Trails. Both datasets were split into sentences and processed to eliminate boilerplate HTML or headers. All sentences containing a single quantity were incorporated into the collection. The entire test data is manually created and tagged. In the following, we describe the query formulation and annotation task.

Query formulation: Given access to the concept/unit index and the value distributions, annotators were tasked to formulate quantity-centric queries. They were instructed to scan the entire index for possible synonyms to a given concept when formulating a query and keep track of the synonyms in a list. For example, if one chooses "Microsoft Surface Earbuds" with the unit "pound sterling", the annotator scans the other concepts inside the concept/unit index with "pound sterling" as unit to detect relevant ones, e.g., "Earbuds", "Microsoft headphones". In the subsequent stage, the value distributions of all selected concepts are consolidated into one and presented to the annotator. The annotator is then instructed to choose three values for equal, less than, and greater than queries, in such a way that supporting sentences for the query is present within the value distribution. In the final stage, the annotator will formulate the query in natural text, e.g., "Microsoft Surface Earbuds lower than 179 pound sterling". The annotators have access to the dictionary of surface forms for units and conditions to help query formulation.

Candiate list generation : For each query, a list of candidates relevant sentences was generated using the *concept/unit index*. All sentences related to

the concept and its synonyms were filtered based 1154 on the query value and condition. The filtering 1155 is done automatically based on the query value 1156 and numerical condition to lower the effort of 1157 annotation. We recognize that the quality of the 1158 candidate set relies directly on how effectively the 1159 quantity extractor captures associations between 1160 quantity and concepts. We observed that although 1161 the extractions for financial data were of high 1162 quality, in the medical domain, several quantities 1163 were overlooked. In both datasets, there is no 1164 guarantee that the candidate list is comprehensive 1165 and covers all relevant instances. 1166

> Annotation: An annotation guideline was devised for a consistent annotation of ambiguous cases and is published with the dataset. Annotators were presented with a list of candidate sentences for each query and were tasked to mark the relevant sentences. The marked sentences are used at ground truth for subsequent evaluation.

C.2 Semantic and Lexical Queries

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

1179

1180

1181

1182

1183

1184

1185

1186

1187

1188

1189

The queries are categorized into four types: *seen*, *unseen*, *expansion*, and *w/o surface form*. The lexical queries fall under the categories of *seen* and *unseen*. For such cases, during query formulation, the annotators picked concepts from the *concept/unit index* without the change in surface form. The concepts from the *unseen* category, were removed from the index for data generation and training of the joint neural models. Therefore, it contains lexical queries that were not seen during training. For example, "YouTube channel" is a concept in the *unseen* subset, which means all instances of "YouTube channel" were removed from the *concept/unit index* before data generation.

Semantic queries fall under the categories of expan-1190 sion and w/o surface form and were slightly harder 1191 to formulate, thereby, fewer instances of them are 1192 present in the data. For expansion queries, a con-1193 cept from the lexical set was chosen to expand to 1194 one of its supersets or synonyms. For example, 1195 "social media channel" is a semantic concept from 1196 "YouTube channel". These expansions were used to 1197 formulate queries that did not have a lexical match 1198 in the database and often included a superset of 1199 1200 many concepts. In the case of "social media channels", the model should be able to retrieve other 1201 social media channels like "Facebook" as well as 1202 "YouTube". In the case of lexical models based on BM25, the difference is evident in Figure1a, 1204

where the models show great performance on seen1205and unseen subset, but if the same queries are converted to their semantic counterpart, as in expansion, the models fail to retrieve the correct result.1206W/o surface form are other semantic queries that1207were formulated independent of the lexical queries.1209

1211

1212

1213

1214

1215

1216

1217

1218

1219

1220

1221

1222

1223

1224

1225

1226

1227

1228

1229

1230

1231

1232

1233

1234

1235

1236

1237

1238

1239

1240

1241

C.3 Implementation

The code is implemented in Python 3.10.9 and PyTroch 1.13.1. The general sentence splitting and text cleaning were performed with SpaCy 3.6⁴. As mentioned before we use the CQE library ⁵ for quantity extraction. Evaluation and metrics were computed with the help of pytrec_eval library (Van Gysel and de Rijke, 2018). In the following, we discuss the implementation details for each model separately.

BM25 models: We use the Okapi BM25 package ⁶ for all BM25 variants. The QBM25 and BM25_{*filter*} are variations of Okapi BM25 designed to include a numerical index for ranking and filtering. The parameters of BM25 were tuned to each of our datasets separately, as presented in Table 4. The latency values are computed with plug-ins for an opensearch ⁷ instance on a desktop computer with 15GB of RAM. In comparison to the dense models, the lexical models do not require specific hardware architectures.

Table 4: Hyper parameters of BM25-based models on the benchmark datasets.

| | FinQuant | MedQuant |
|-----------------|--------------------|--------------------|
| BM25 | b = 0.5, k1 = 0.5, | b = 0.5, k1 = 0.5 |
| $BM25_{filter}$ | b = 0.75, k1 = 1.5 | b = 0.75, k1 = 1.5 |
| QBM25 | b = 0.5, k1 = 0.5 | b = 0.5, k1 = 0.75 |

Cohere baseline: We used the Cohere API ⁸ for Cohere $_{v3}$ embeddings. Query embeddings were used to encode the queries and the document embeddings to encode the collection.

ColBERT models: (Khattab and Zaharia, 2020) supplied the trained checkpoint for the base ColBERT model. For fine-tuning on augmented data, the model was initialized with this base checkpoint. The checkpoint was employed for

⁴https://spacy.io/ DLA: 11.02.2024

⁵https://github.com/vivkaz/CQE DLA: 11.02.2024

⁶https://pypi.org/project/rank-bm25/ DLA: 11.02.2024

⁷https://opensearch.org/ DLA: 11.02.2024 ⁸https://cohere.com/ DLA: 11.02.2024

the evaluation of both ColBERT and QColBERT. 1242 $ColBERT_{ft}$ was fine-tunned using the training 1243 script the from the official repository 9 . The 1244 code in the repository was modified to establish 1245 an endpoint for QColBERT incorporating a quantity index. We did not perform extensive 1247 hyperparameter tuning except for the learning rate 1248 and mainly used the parameters advised by the 1249 authors for both FinQuant and MedQuant datasets. 1250 We fine-tuned the joint $ColBERT_{ft}$ for 2 epochs, 1251 with a batch size of 256 and a learning rate of 1252 1e-05 on a server with four A-100 GPUs and 40GB 1253 of memory. The evaluation and benchmarking 1254 for latency were performed on the same server, 1255 utilizing all four GPUs. 1256

1257

1258

1259

1261

1263

1264

1267

1268

1269

1274

1275

1276

1277

1278

1279

1280

1281

1283

1284

1285

1287

1288

SPLADE models: SPLADE_{*ft*} was also fine-tuned using the training script by the authors ¹⁰. The pre-trained checkpoint was acquired from Hugging Face ¹¹ and utilized for both the SPLADE model and QSPLADE. Scripts from the official repository were adjusted to add numerical index for QSPLADE. Similar to ColBERT, we conducted limited hyperparameter tuning, mainly focusing on the learning rate. We fine-tuned SPLADE_{*ft*} for 2 epochs using a batch size of 240, a learning rate of 2e-5, and a weight decay of 0.01. The fine-tuning was conducted on a server with four A-100 GPUs and 40GB of memory. The evaluation and benchmarking for latency were performed on the same server, utilizing all four GPUs.

For all disjoint rankers, QBM25, QColBERT, and QSPLADE the quantity impact parameter of α is set to 1, such that the impact of term and quantity ranking are equal.

Generated data: Based on the combination of augmentation methods the size of training data would vary. In all cases, we saved a small sample of 1000 queries for validation. There were 40,732 and 20,376 concept and unit pairs considered for query generation in FinQuant and MedQuant, respectively. If concept expansion is applied these numbers would double to account for queries on expanded concepts. We set the sample size n to 2, meaning that for each query two positive and neg-

⁹https://github.com/stanford-futuredata/ ColBERT DLA:11.02.2024

¹⁰https://github.com/naver/splade DLA:11.02.2024 ¹¹https://huggingface.co/naver/



Figure 7: The effect of task-specific fine tuning on models attention to quantity tokens. In the masked variants either the unit or the value of the sentences in the collection is masked.

ative samples were chosen from the data without augmentation. As a result, based on augmentation methods, additional n = 2 samples would be added for unit and value permutation, a total 3n per query.

1289

1291

1293

1294

1296

1297

1298

1299

1301

1302

1304

1305

1306

1307

1309

1310

1312

1313

1314

1315

1317

1318

C.4 Effect of Fine-tuning

To assess the impact of task-specific fine-tuning on the internal ranking strategy of the dense models, we evaluate two masked versions of the data.

Mask value: In this scenario, we mask all values in the collection with the [MASK] token before running the evaluation. This task aims to determine the extent to which the model depends on the value token for retrieving the correct sentence.

Mask unit: Here, we mask unit tokens in the collection before running the evaluation with [MASK] token. This task is intended to observe the impact of unit comparison on the final ranking.

We compare the base version of the dense models with their fine-tuned version on the different masking of the FinQuant dataset. The results for the ColBERT models are shown in Figure 7a and for SPLADE models in Figure 7b. In both cases, the fine-tuned version exhibits a more significant drop in performance compared to the base models. This indicates that after fine-tuning, the model becomes more dependent on the quantity of tokens, values, and units, in the text to identify the relevant sentence.

C.5 Ablation Study on Augmentation Methods

To check the effect of different augmentation strate-
gies, we perform an ablation study, by fine-tuning13191320

 $^{{\}tt splade-cocondenser-ensembledistil DLA:} 11.02.2024$



Figure 8: Ablation study on different augmentation methods, where *value* and *unit*, refer to value and unit permutation and *concept* refers to concept expansion.

the neural models on data generated using a com-1321 bination of different strategies. The main points of 1322 variabilities are concept expansion in the query gen-1323 eration process and value and unit permutation for 1324 sample generation. The results for $ColBERT_{ft}$ and 1325 $SPLADE_{ft}$ on FinQuant dataset is demonstrated in 1326 Figures 8a and 8b, respectively. no perturbation 1327 refers to the case where no data augmentation was 1328 applied and only the positive and negative samples 1329 from the original sampling are used for training. An interesting trend is the detrimental effect of 1331 value permutation. The value permutation on 1332 its own enhances the performance of the base 1333 model. However, as soon as it is accompanied by 1334 1335 other augmentation methods the performance degrades slightly. The best combination for both the 1336 SPLADE and ColBERT model is unit permutation 1337 and concept expansion, both of these augmentation 1338 techniques on their own also provide a larger boost 1339 1340 in comparison to value permutation. To this end, the variant of the models presented for evaluation 1341 as $ColBERT_{ft}$ and $SPLADE_{ft}$ are trained on unit 1342 permutation and concept expansion subset. We 1343 find this behaviour rather surprising and counter-1344 intuitive. Usually, the performance of neural mod-1345 els increases with the amount of data presented for 1346 a given task, however, perturbing the values does 1347 not seem to enhance the performance as expected. 1348 This can be related to the internal representation 1349 of the neural models, which is hindering their abil-1350 ity to correctly learn quantity semantics. In future 1351 work, we aim to test the effect of dedicated numer-1352 ical embedding and language models for this task. 1353 1354