A REASONING-BASED APPROACH TO CRYPTIC CROSSWORD CLUE SOLVING

Anonymous authors

004

005

010 011

012

013

014

015

016

017

018

019

021

023 024

025

035

041

048

Paper under double-blind review

ABSTRACT

Cryptic crossword clues are challenging language tasks for which new test sets are released daily by major newspapers on a global basis. Each cryptic clue contains both the definition of the answer to be placed in the crossword grid (in common with regular crosswords), and 'wordplay' that *proves* that the answer is correct (i.e. a human solver can be confident that an answer is correct without needing crossing words as confirmation). This work describes an LLM-based reasoning system built from open-licensed components that solves cryptic clues by (i) hypothesising answers; (ii) proposing wordplay explanations; and (iii) using a verifier system that operates on codified reasoning steps. Overall, this system establishes a new state-of-the-art performance on the challenging Cryptonite dataset of clues from The Times and The Telegraph newspapers in the UK. Because each proved solution is expressed in Python, interpretable wordplay reasoning for proven answers is available for inspection.

1 INTRODUCTION

Recent advances in computational models have significantly improved their ability to handle diverse natural language tasks involving complex syntactic and semantic interpretations. Despite these strides, machines continue to fall short of human performance in several areas, most notable being those involving reasoning. This work tackles the relatively under-studied reasoning task of cryptic crossword solving, which is a popular activity across the world, with multiple papers in the UK, Australia, India and elsewhere featuring daily puzzles for readers to solve.

The following is an example of a cryptic clue¹ of *moderate* complexity: $\frac{1}{1000}$

034 clue (4D): "Cut up over politician on the French case (7)"

O36 Solvers must interpret the clue to understand the supporting wordplay to arrive at the answer from two directions (see also Figure 2a for a visual depiction of the reasoning involved):

038 definition: "Cut up over politician on the French {case}"
039 wordplay: "(AXE) < (cut, <up), MP, LE (the, in French)"
040 answer: "EXAMPLE"</pre>

In the above, the reasoning steps are: (i) identifying which part of the clue is the definition (highlighted with curly braces) that acts as a regular crossword clue; (ii) parsing the remainder of the clue for the wordplay - here, for instance, (a) the indicator word 'up' tells us to reverse the word resulting from 'Cut' (since this is a down-oriented clue), (b) this is 'over' the abbreviation of Member of Parliament 'MP' (a politician in the UK), (c) and 'on' a French translation of 'the' (common knowledge); and (iii) finally assembling the parts to 'prove' that the correct answer is 'EXAMPLE' (this agrees with the definition span, with the correct number of letters).

In this work, we take our cue from the effectiveness of provers coupled with verifiers for mathematical reasoning tasks (Jiang et al., 2023). We tackle the cryptic crossword clue solving task using an LLM to (i) suggest answer candidates; (ii) create informal proofs (i.e. coming up with wordplay); and (iii) perform a formalisation process (which rewrites the wordplay logic in Python). The proposed solutions are then checked with a verifier for validity.

¹The Times UK (6-March-2024), Cryptic #28857 clue 4D



• Local models for cryptic clue tasks - We show how local LLMs can be fine-tuned to produce answer candidates, and wordplay suggestions, and then prompted to perform Wordplay formalisation. Following an approach akin to mathematical statement formalisation, but where there are *less than 10* examples of 'good proofs' available, our novel pipeline was specifically engineered to avoid 'reasoning steps' becoming stuck in dead ends.

• **Python domain-specific verifier** - Using the output of the formaliser, the verifier presented here deconstructs the Python AST, so that it can evaluate each assert statement on a line-by-line basis. We believe that this is somewhat novel, since it enables the verifier to not only indicate whether the proof is valid overall, but also point to specific failures (used to regenerate failed formalisations) on all proof lines simultaneously.

To promote further study in this area, all code for the formaliser and domain-specific verifier is made publicly available.

094 095 096

098

100

085

087

090

092

093

2 RELATED WORK

099 2.1 REGULAR CROSSWORDS

Non-cryptic ("regular") crosswords are known throughout the world, and are the predominant type found in newspapers in the U.S.A. One key difference from cryptic crosswords is that individual regular crossword clues are generally not 'standalone' - there may be a number of different answers that fit the given clue. The key to solving regular crosswords is thus the interaction between answers (i.e. the crossing-words), which allows for planning/backtracking to enable solving rates in the high 90% range (Wallace et al., 2022).

107 This work, in contrast, focuses on the solving of clues on a standalone basis, which requires elements of reasoning through the wordplay present in cryptic clues.



to those from the Guardian newspaper, and Connor (2024) notes in the Guardian's own blog "The
Times hosts an annual crossword-solving competition and it remains, until such time as the Guardian
has its own version, the gold standard." Moreover, the smaller number (142,000) of clues the dataset
contains have no orientation markings ('across/down'), which are required to make sense of some
wordplay.

For a more in-depth discussion of the decision to focus on the Cryptonite dataset (and not perform testing on the Guardian dataset), please see Appendix A.3. In summary, while the 'Init' split presented in Rozner et al. (2021) has attractive properties (explored there, and in other works), this work specifically targets the wordplay reasoning side of cryptic clues, and since that involves finetuning models on Cryptonite (including Wordplay examples with carefully matched train/val/test splits), this precludes us from doing the same kind of multi-dataset comparisons found elsewhere.

156 2.2.1 RULE-BASED SOLVERS

155

Williams & Woodhead (1979) is an early example of attempting to devise a formal language for describing cryptic clues. However, the linguistic elements of the clues tend to thwart a strictly formal approach.

A more flexible rule-based solver with a manually-crafted probabilistic grammar was introduced in (Deits, 2015; 2022). Building on the assumption that a clue can usually be split into wordplay

"num": 16, "ad": "D",

"answer": "PROPOSAL"

165 166

> 167 168 169

}

162

163

164

Figure 3: An example from the Wordplay dataset (in this wordplay, () * is an anagram indicator). This clue's solution is diagrammed in Figure 2b

"publication": "FT", "setter": "falcon", "author": "teacow",

"clue": "{Offer} of support also broadcast", "pattern": "8",

"wordplay": "PROP (support) + (ALSO) * (*broadcast)",

171 172

170

and definition, the solver tries to find the most probable parse such that the wordplay yields a semantically-similar result to the definition. The logical form of this DSL approach is very appealing. However, it appears limited to solving clues where the wordplay is somewhat simple (due to the combinatorial explosion of possibilities created by longer/more complex clues).

The goal of this work is to use the flexibility of LLMs to enable a far wider range of clues to be attempted, with the aid of a formaliser/verifier to check the solutions.

- 179 180
- 180 2.2.2 LLM-BASED SOLVERS

Cryptonite is a challenging task for LLMs : Efrat et al. (2021) reported that fine-tuning T5-Large (a 770M encoder-decoder model) on Cryptonite's 470k cryptic clue training set achieved only 7.6% test set accuracy, slightly below the 8.6% accuracy of the rule-based clue solver of Deits (2022). Interestingly, prior to 2024, even large-scale Language Models scored very poorly on cryptic clues, likely due to (i) the misleading surface reading of the clues; (ii) the obliqueness of the definitions; and (iii) the reasoning steps required to *prove* the answer correct based on the wordplay.

Recent works, such as Sadallah et al. (2024) and Saha et al. (2024), tackle cryptic crosswords with
more up-to-date local models, and also commercial LLMs. Saha et al. (2024) reports results with
5- and 10-Shot prompting (without fine-tuning models), but also includes a wide-ranging study of
the capabilities of models for crosswords in general. We include experiments that bring the relevant
baselines up-to-date, and also touch on their illuminating Partial Correctness Metrics (which are
relevant when attempting full grids, but not the main focus here).

In this work, we use a pipeline of 9B-scale LMs to produce answer candidates and wordplay suggestions, followed by a third LM to formalise each proposed solution using Python code and then rewrite/update the solutions based on feedback from a purpose-built verifier. In our results, we focus on the 'pure' Cryptonite benchmark: Accuracy is judged based on a Top-1 basis (with the model's single answer being marked wholly correct or not), with no crossing letters being given. Framed as a reasoning task, if the model 'understands' the cryptic solution properly, the answer will be wholly correct - there should be no partial marks.

201 2.3 CODE & REASONING

To compensate for LLM *approximate generation* of logical reasoning, techniques like PAL (Gao et al., 2023) exploit LLMs' facility for writing code to create verifiable reasoning chains. An important influence on this work was also the Draft, Sketch, and Prove framework (Jiang et al., 2023) which uses an LLM to draft and create proofs that are then verified formally.

Informed by the evolution from AlphaCode (Li et al., 2022), in which huge numbers of programs are generated and filtered in order to generate a valid solution, to AlphaCodium (Ridnik et al., 2024), in which solutions are iterated upon and involving much less computation, this work uses a verifier that can feed back 'hints' to the formalising LLM, so that the task of re-writing nearly-valid proofs is made easier.

212

- 213 2.4 WORDPLAY DATASET
- The Wordplay dataset (Andrews, 2024) an example from which is given in Figure 3 consists of data gathered from websites where cryptic crossword enthusiasts post solutions on a daily basis for

each of the major publications. Each completed puzzle is annotated by an individual, identifiable author/solver that lists the approximately 30 clues with their definition, wordplay and answer
fields. Note that the authors each chose their own 'standard' for writing out the wordplay, leading
to a significant variation in wordplay annotation style between solvers : Even though the underlying reasoning is clear, solvers do not annotate in the same style (even across time). The Wordplay
dataset deliberately follows the train, validation, and test splits defined by Cryptonite.

222 223

224

3 Methods

The overall system described in the work is illustrated in Figure 1. The order of operations for the pipeline was chosen based on watching human solvers - who report going through the following steps: (a) attempting to parse the clue in a number of ways, trying to isolate the definition from the wordplay; (b) seeing which parts of the wordplay they are most confident about; (c) 'having a hunch' of the final answer; and (d) gaining a full understanding of how a clue's wordplay works (such that the function of every element can be explained) as proof of the overall process.

Observations of the behaviour of GPT-4 using Chain-of-Thought prompts Wei et al. (2023) (or more recently 'o1'), suggest that even very capable models tend to fixate early on during the reasoning process, and are only rarely able of completely re-hypothesising. Theses LLMs also frequently becomes caught up with the literal ('surface') meaning of the clue, which is often misleading. The combination of these elements caused us to re-consider how to approach this kind of problem.

By organising our system's pipeline to hypothesis candidate answers as the first step (so that the models must try to fit the reasoning to the answer, with varying degrees of success) bakes rehypothesisation into the process. Another aspect that required care was the looseness of the language used in cryptic clues, which may stop even valid 'reasoning' from being provable.

240 241

242 3.1 CANDIDATE answer GENERATION

Our first step to solving a given cryptic clue is to generate multiple answer candidates from the original clue, pattern and ad (across/down) fields. For this task, we fine-tuned a Gemma2 9B base model (Gemma Team & Google DeepMind, 2024) using the LoRA (Hu et al., 2021) implementation provided by the unsloth package (unsloth.ai, 2024). The model was trained for 1 epoch on the Cryptonite training set of approximately 470,000 examples.

For each clue being evaluated, we generate 20 valid answer candidates, where candidates that did not match the pattern were immediately rejected and regenerated, and those not contained in the crossword words list (Beresford, 2000) were filtered out². The number of candidates was chosen to balance generation cost with likelihood of the correct answer appearing in the candidate list - see Figure 7 for a cumulative frequency analysis. The list of candidates was then grouped so that the frequency of each answer could be found - enabling statistics to be collected.

3.2 GENERATION OF definition AND wordplay SUGGESTIONS

To train the wordplay suggestion model, which translates each answer candidate into multiple definition and wordplay suggestions, we make use of the Wordplay dataset of Andrews (2024). For this task, we fine-tuned another Gemma2 9B base model using LoRA. The model was trained on 4 epochs on a set of approximately 16,800 examples (consisting of solution explanations of puzzles from The Times and The Financial Times from selected authors in the Wordplay dataset).

3.3 PYTHON FORMALISATION

Rather than create a dataset with many examples of formalisation, here we use in-context prompting with less than 10 examples of the formalisation style required. In preliminary work, we concluded that the available Gemini-Flash LLM was not capable of using a (novel) cryptic crossword domain specific language ("DSL") through in-context learning with so few examples. In contrast, we found

267 268 269

254 255

256

257

258

259

260

261 262

263 264

265

²This rejection of invalid words here is not 'cheating' since we do not use the dictionary to suggest words, rather it is only used to weed out actively proposed non-words from a short-list.

```
270
                                                                 Gemma2 answer candidate
                def proof(answer="DELVE",
271
                          clue="research done, primarily on most of magical beings",
272
                          pattern='5'):
273
274
                                                                       Gemma2 wordplay
                  ....
275
                  definition: research done, primarily, on most of magical beings
276
                  wordplay: D[one] (primarily) ELVE[s] (magical beings, most of)
277
                  .....
278
                                                                        LLM formalisation
279
                  assert action_type("primarily", Action.INITIALS)
280
                  assert "DONE"[:1] == "D"
281
                  assert is_synonym("magical beings", "ELVES")
                  assert action_type("most of", Action.REMOVE_LAST)
282
                  assert "ELVES"[:-1] == "ELVE"
                  assert "D"+"ELVE" == "DELVE"
284
                  assert is_synonym("research", "DELVE", pattern='5')
                proof()
287
              Figure 4: Python proving: answer candidate \rightarrow wordplay \rightarrow LLM formalisation
289
       Action=Enum('Action',
290
                      'ANAGRAM, REMOVE_FIRST, INITIALS, REMOVE_LAST, '+
291
                      'GOES_INSIDE, GOES_OUTSIDE, REVERSE, SUBSTRING, HOMOPHONE')
       # External definitions
292
       def is_synonym(phrase:str, test_synonym:str, pattern:str='') -> bool:
293
          # True if 'test_synonym' is a reasonable synonym for 'phrase',
294
          # with letters optionally matching 'pattern'
295
       def is_abbreviation(phrase:str, test_abbreviation:str) -> bool:
296
          # Determines whether 'test_abbreviation' is
          # a valid abbreviation or short form for 'phrase'
297
       def action_type(phrase:str, action:Action) -> bool:
298
          # Determines whether 'phrase' might signify the given 'action'
299
       def is_anagram(letters:str, word:str) -> bool:
300
          # True if 'word' can be formed from 'letters' (i.e. an anagram)
301
       def is_homophone(phrase:str, test_homophone:str) -> bool:
          # Determines whether 'test_homophone' sounds like 'phrase'
302
303
                       Figure 5: External functions available via In-Context Learning
305
306
307
       that the LLM could be prompted to produce Python code with relative ease, so the approach taken
308
       was to frame a declarative-style-DSL as Python function calls within assert statements. The LLM
       was found to be able to reliably produce syntactically correct Python, and use the 'external functions'
309
       that had been described (as illustrated in Listing 5) to form logical sequences of declarations, which
310
       could then be parsed line-by-line by manipulating the Python abstract syntax tree ("AST"). An
311
       example of the Python DSL being generated by the formalisation LLM is given in Figure 4, with the
312
```

- workings of the clue solution being illustrated in Figure 2c.
- To formalise wordplay into Python 'proofs' of the correctness of solutions, we used Google's Gemini-Flash-1.5-001 LLM (a pinned model version) during development. This model was initially chosen instead of a frontier-tier model since the formalisation task should not require much inventiveness/reasoning: the actual required steps are already present in the wordplay, the task is *merely* to translate to Python. To determine whether the choice of Gemini-Flash was a limiting factor, we subsequently tested an unmodified Gemma2-9B-it model on the same task.
- In terms of the DSL itself, the back-end to the is_synonym and is_homophone functions consists of calls to simple language models. The action_type function performs a nearest-neighbour match against list of indicator words, and the is_abbreviation function performs a look-up against a list of abbreviations - both sourced from Deits (2022). For string manipulation actions (such as 'REVERSE'), the LLM formaliser itself was capable of producing correct output unaided.

324 AssertionError: assert: is_abbreviation('an Artist', 'RA') : \ 'an Artist' does not have a valid abbreviation; \backslash 325 'RA' is an abbreviation for : $\$ 326 artist, artillery, Royal Artillery, gunners, painter 327 AssertionError: assert action_type('goes crazy', Action.ANAGRAM) : \ 328 'goes crazy' does not suggest Action.ANAGRAM, but 'crazy' does AssertionError: assert action_type('worked', Action.HOMOPHONE) : \ 'worked' does not suggest Action.HOMOPHONE, but may be Action.ANAGRAM 330 331 332 Figure 6: Illustrative AssertionError responses (with hinting) from the verifier 333 334 3.4 IN-CONTEXT LEARNING 335 336 To produce Python code that could be sent to the prover, the LLM was prompted in an In-Context 337 Learning ("ICL") manner. This consisted of the following parts: 338 1. Cryptic crossword rubric to explain to the LLM what the principles were behind the fields 339 such as clue, definition, wordplay, etc. 340 2. Many-shot clue \rightarrow wordplay examples (20-30 given) 341 342 3. The 'external functions' rubric shown in Figure 5 343 4. Few-shot wordplay \rightarrow Python formalisations (6 examples given) 344 5. The actual clue, answer, definition and wordplay being formalised 345 Gemini-Flash did not appear to be particularly sensitive to the prompting style used, except in 346 the 'handover step' (between problem description and model generation) where several trials were 347 needed to obtain the final function definition in the required format consistently. Further details of 348 all the ICL prompts are given in Appendix A.5. For the final Gemma2-9B-it formalisation runs, the 349 same prompts were used unchanged (with no other tuning/training). 350 351 3.5 PROOF VERIFICATION WITH HINTING 352 353 The verifier implemented here must decide whether a given formalisation is valid, and report any 354 errors found to iteratively improve the Python code as feedback to the LLM formaliser in a cycle, as seen in Self-Debug (Chen et al., 2023), and AlphaCodium (Ridnik et al., 2024). Examples of 355 assertion failures, with constructive hinting, are shown in Figure 6. 356 357 This cycle is repeated until a formalisation is validated (zero assertion failures, considered a 'SUC-358 CESS' with the answer having been proved), or max_rewrites=2 is reached. If no Python 359 formalisation can be validated, then the fallback answer is used (defined as being the most fre-360 quent answer amongst the original candidates produced in the first stage of solving the clue). 361 362 3.6 PARTIAL CORRECTNESS METRICS One interesting direction explored in Saha et al. (2024) was the performance of LLMs on cryptic 364 clues if some of the letters were known (as would be the case if an entire grid were being solved). The conditions examined were with 25%, 50% and 70% of letters 'known'. Based on observa-366 tions working with the proposed system for single clues (and more general experience of crossword 367 solving), we approached this problem in two ways. 368 For the 25% level of letters 'known', it was a simple matter to use our existing system with candidate 369 answers which didn't match the known letters filtered out. For the higher levels of known letters, we 370 instead used the FastText embedding method of Mikolov et al. (2018) to find the nearest neighbour 371 answer within The UK Advanced Cryptics Dictionary, Beresford (2000), by comparing against 372 the embedding of the raw clue phrase itself. 373 374 The 'Partial Correctness' results, while not being a core thrust of our reasoning approach but interesting in their own right, are given in Appendix A.5.7. 375



Figure 7: Statistics of answer candidate list, as more candidates generated

4 EXPERIMENTS

4.1 GEMMA2 9B answer CANDIDATE GENERATION

During the initial experimental phases of fine-tuning local models for the answer generation task it was discovered that -base models scored more highly than -it models. This might be explained by observing that instruction fine-tuning may (to some extent) penalise off-the-wall answers, which may be essential for our task. In addition, we also observed that while the Top-1 candidate from a model generating with a temperature t = 0.5 had high accuracy, it was beneficial to run candidate generation with t = 1.0 (even though the Top-1 accuracy was lower in this case) - since having a wider spread of answer candidates was useful for our pipeline overall.

Figure 7a shows that the probability of the gold answer being among the candidates produced is (unsurprisingly) monotonically increasing in the number of independent samples. It also shows that this process is not yet asymptotically limited, although slowing down with increasing n.

Figure 7b shows that choosing the highest-frequency answer candidate can be a very effective strategy. However, there is a clear limit to this idea: There is a significant probability that cryptic crossword answers are in the long tail of potential answers. Indeed, intentionally creating misleading clue 'surfaces readings' is a hallmark of good cryptic clue setting.

411 412 413

391 392 393

394

396

4.2 GEMMA2 9B wordplay CANDIDATE GENERATION

Since wordplay is so flexible, it is difficult to evaluate it for accuracy against other examples
(without, say, a large LLM to evaluate the differences). However, good wordplay should result in
good formalisations, so evaluation is available on an end-to-end basis.

One key assumption in the system proposed here is that a correct answer should lead to interpretable wordplay, whereas an incorrect answer candidate should give rise to unformalisable/unverifiable wordplay. The following typical example illustrates how the correct answer leads readily to correct wordplay (the workings of this clue are illustrated in Figure 2d), whereas trials with an incorrect answer candidate (which was, in fact, the most frequent candidate for this clue) give clearly unverifiable wordplay:

```
423
      clue: "wader woman has on (5)"
424
          definition: "{wader} woman has on" # Same for all
425
        answer: "HERON"
                          # correct answer
426
          wordplay: "woman (HER) has on (ON)"
427
428
                          # incorrect answer - 3 trials shown
        answer: EGRET
429
          wordplay: "woman (HER) on top of (REG - another word for on, \
430
                      as in 'do you have the heating on?')"
          wordplay: "EG (woman has) + RET (on)"
431
          wordplay: "woman (HER) has on/around (EG) - a wader bird"
```

432

450

451 452 453

454 455

456

457

458

459

460 461

462

463

464

465 466

467 468

469

470

			Validation		Test		
Model	samples	Overall	Quick	Hard	Overall	Quick	Haro
Rule-based (*) T5-large (770M) FT (*)	26k 26k	8.3% 7.4%			8.6% 7.6%	13.5% 12.8%	5.8° 3.4°
Gemma2-9B-it 5-shot	1000	5.7%	11.5%	5.2%	4.5%	10.5%	4.0
Gemini-Flash 5-shot	1000	6.6%	12.5%	6.1%	6.5%	11.8%	6.1
GPT-40 5-shot	1000	29.8%	45.0%	28.5%	27.6%	47.4%	26.0
Gemma2-9B FT	1000	21.7%	28.8%	21.1%	15.9%	38.2%	14.1
Gemma2-9B freq (#=20)	1000	26.6%	31.3%	26.2%	25.5%	55.3%	23.1
(AB) logprob answer	500	23.9%	35.9%	22.9%	22.7%	55.3%	20.1
(AB) logprob wordplay	200	21.0%	15.4%	21.4%	20.5%	46.7%	18.4
Gemini-Flash Formaliser	200	28.0%	23.1%	28.3%	32.5%	46.7%	31.4 27.6
Gemma2 9B-it Formaliser	200	26.0%	23.1%	26.2%	29.0%	46.7%	

Rows (*) are as reported in Efrat et al. (2021); The Hard columns are for the non-Quick clues

4.3 CRYPTONITE RESULTS (TOP-1 EXACT MATCH)

In this work, we focus our testing on using the Cryptonite dataset of Efrat et al. (2021) as a benchmark, with the Top-1 exact-match results shown in Table 1. As in Saha et al. (2024), due to computational constraints, we performed sampling of the validation and test sets, using fewer than the full 26k examples available. One standard deviation at 1000 samples is $\approx \pm 1.5\%$, and at 200, $\approx \pm 3.3\%$.

The 5-Shot results in Table 1 show that:

- GPT-40 (2024-11-25) gives stronger results than those of GPT-4-Turbo (2024-04-09) given in Saha et al. (2024) so this is an updated baseline;
- The updated GPT-40 results show surprisingly strong performance on the validation split (unfortunately, the composition of this commercial model's training data is unknown);
- Gemini-1.5-Flash-001 (which was used in development of the formaliser) is not particularly good at solving cryptic clues in itself;
- The Gemma2-9B model gets a large uplift from fine-tuning on the Cryptonite training set (compare the 5-Shot figures to the later Gemma2-9B FT ones).

The Gemma2-9B FT accuracy figures are for the first result returned by the fine-tuned Gemma2 model. In contrast, the Gemma2-9B freq accuracy figures are for the most common (i.e. highest frequency) result among the Gemma2 answer candidates (for which 20 samples were generated for each clue). These voting-based results show impressive performance, which would have exceeded prior state-of-the-art results for open-licensed models on their own.

- Going beyond single models, the Gemini-Flash Formaliser demonstrates Top-1 exactmatch performance of 32.5% for the Cryptonite Test set, establishing a new state-of-the art result, even against the updated baselines. Moreover, the results of the non-fine-tuned Gemma2-9B-it Formaliser also beat the previous state-of-the-art results - which is perhaps an even stronger statement about the capabilities the system described here, since in this case Gemma2-9B models have been used throughout the solving process, showing that it is be possible to achieve very competitive cryptic crossword solving results through reasoning with off-line, open-licensed models.
- The formaliser results are (surprisingly) relatively worse for Quick clues. This seems to be related
 to the fact that the agreement/frequency-based Gemma2 freq model is very strong on these clues,
 and any 'contribution' from the formalising/verification procedure is likely to overrule a good baseline result, due to erroneous verification of 'proofs' that are not valid.

486 4.4 ABLATIONS

The lines in Table 1 marked '(AB)' are ablations. Both utilise the measurement of average *logprob* of the output tokens given by the relevant model.

The first ('logprob answer') shows the results of using the Gemma2-9B FT model from above, with the candidate answer being chosen from the list of 20 possibilities according to highest *logprob*. Since answers are typically very short, this method is similar to the frequency-based selection model.

The second ('logprob wordplay') shows the results of evaluating the Gemma2-9B FT model that generates wordplay hypotheses, and choosing an answer based on the highest *logprob* according that generating model. Somewhat unexpectedly, this was not as effective as might be assumed from the generated wordplay seen in Section 4.2 - where the wordplay for wrong answers looks absurd. Examining samples of the wordplay most favoured by pure *logprob* order, it seems that the generating LLM finds simply-worded but *completely fictitious* wordplay quite likely.

Both of these ablations demonstrate that the formalisation and verification steps are essential components in our system - we cannot shortcut them using a 'dumb ranker' in the pipeline.

502

504

505

506 507

508

509

510

511

512 513

514

515

4.5 OBSERVED LIMITATIONS OF THE SYSTEM

The verifier implemented for this work does not detect a number of potential errors in the Python function it analyses:

- The entire Python function might consist of comments : Nothing triggers assert
- The Python function contains conditional execution, routing around assert statements
- Occasionally, the hint assert XYZ failed results in the re-write : assert XYZ==False
- The proof may be logically disconnected, with left-hand-side terms not being supported / justified by right-hand-side terms in other lines of the code

These issues do not appear insurmountable, given time and effort. It should be noted that since the formalising LLM is only being used In-Context there is little chance that the above issues are being systematically abused (which would almost certainly happen if there was learning-in-the-loop in a Reinforcement Learning setting).

516 517 518

519

5 CONCLUSIONS

It is hypothesised that the next-token-prediction task may be insufficient to get machines to reason and plan (Kambhampati, 2024). Our choice of platform for reasoning experiments was that of cryptic crossword clue solving - and we have demonstrated early success with a *system* that include both LLMs, verifiers and coding aids.

We believe that our results validate our overall approach to codification of the cryptic crossword problem domain : Generating answer candidates and wordplay suggestions followed by production of code via an LLM-based formalisation process, verification using Python code analysis tools, and iterative prompting of the LLM formaliser proved quite effective.

We were happy to discover that our development work using the Gemini-Flash LLM as a formaliser
was directly transferable to a non-fine-tuned version of the open-licensed Gemma2-it model for the
same role, with little loss of performance, enabling the whole pipeline to be run locally. The weakest
link in the chain was, predictably, getting the 'Aha' of wordplay creation to work - humans can
still generate wordplay that is beyond the capabilities of current models.

The authors sincerely hope that this work sparks interest in the cryptic crossword domain, which presents an array of interesting and challenging reasoning problems on which fruitful research can take place, even for those with limited computation budgets.

- 537
- 538
- 539

540 6 ETHICS STATEMENT

542 6.1 SOCIETAL IMPACT

There are many current cryptic crossword enthusiasts that would potentially not welcome AIenabled solvers to 'take over' their favourite pastime. In particular, when taken further, this line of work would potentially disruptive to public leaderboards that rank people according to the time taken to solve puzzles 100% correctly. However, there is currently little risk of LLM cryptic solvers as being anything more than comic relief for current experts.

6.2 POTENTIAL BIAS IN FAVOUR OF NATIVE ENGLISH SPEAKERS

While the English language has a high capacity for ambiguity and wordplay overall, making this type of crossword possible, cryptic crosswords also exist in other languages (Wikipedia contributors, 2024). In addition, although solving cryptic crossword answers may be very difficult (even for native English speakers), understanding the answer from given wordplay is much simpler.

- 7 Reproducibility
- 7.1 DATASETS AND CODE

The following outline the efforts that have been made to ensure reproducibility:

- **Datasets** Resources such as dictionaries used, and the Cryptonite and Wordplay datasets are available online, via the sources referenced in the main text.
- System Code & LLM Prompts Python code for the complete end-to-end system will be made available on publication. The prompting required for the LLM formalisation (arguably a form of programming) are included in the Appendix.
- **Models used / training** The models referenced in this work are available with open weights (the Gemma2 models), or via API (the Gemini-Flash model, with a pinned version number). The training procedures are outlined in the text, and therefore have some degree of reproducibility. The fine-tuned Gemma2 models will be made available in any case, on publication.

7.2 COMPUTATIONAL REQUIREMENTS

The Gemini LLM was accessed by API, and the total spend to create the results in this paper was
less than \$100 USD, with each prompt round-trip taking around 5 seconds. The Fine-Tuning of
the Gemma2 9B model took around 24 hours for a full Cryptonite training run, and 8 hours for the
Wordplay dataset runs. Thus, the single-GPU model runs totalled less than \$50 USD.

594 REFERENCES 595

Martin cryp	Andrews.	Wordplay	Dataset,	2024.	URL	https	s://g:	ithub	.com/n	ndda/	
J Ross H http	Beresford. Th	e UK Advanc	ed Cryptio m/wordf	cs Diction	ary. Tec	hnical re	eport, p	ublishe	d online	, 2000.	
Xinyun to sel	Chen, Maxw If-debug, 2023	ell Lin, Nath 3. URL http	anael Scha bs://ar	ärli, and D xiv.org	enny Zh [/abs/:	nou. Tea 2304.(iching l 05128.	large laı	iguage i	nodels	
Alan word theo cryp	Connor. remain th guardian.c	Devious e last thin com/cross sword-ai-	humour g AI ca words/c conquer	and pai in't conc rosswor r-human	nful p quer?, rd-blc -solve	uns: 2024. og/202 ers-ar	will UR 4/nov	the c L ht: /04/ cial-:	ryptic tps:// intel]	cross- /www. Ligenc	e-softwar
Robin I 2015.	Deits. rdeits/c	ryptics githu	b code rep	o.http	s://gi	thub.	com/r	deits	/cryp	tics,	
Robin I Cryp	Deits. rdeits/c pticCross	rypticcrossw vords.jl,2	ords.jl git 2022.	hub code	repo. ht	ttps:/	//gitl	hub.c	om/rde	eits/	
Avia Eff for ex Scott Lang 2021. http	rat, Uri Shaha treme ambigu Wen-tau Yih uage Process Association ps://aclar	im, Dan Kilm nity in langua (eds.), <i>Proce</i> ing, pp. 4186 for Computation thology.	nan, and O ge. In Mar <i>edings of</i> –4192, Or ional Ling org/202	mer Levy. ie-Francin the 2021 nline and l guistics. d 1.emnlp	Crypto ne Moens <i>Conferer</i> Punta Ca oi: 10.18	nite: A s, Xuanj <i>ice on E</i> ana, Doi 8653/v1 1.344.	cryptic ing Hua <i>Empirica</i> minicar /2021.e	crosswo ang, Lu <i>al Meth</i> n Repub emnlp-n	ord benc cia Spec ods in N lic, Nov nain.344	chmark ia, and <i>latural</i> vember . URL	
Kathryn novel doi: psyc	n J. Friedland l area of cryp 10.3389/fpsy chology/ar	er and Philip tic crossword g.2016.00567 cticles/1	A. Fine. 1 solving. 7. URL 1 0.3389/	The grou Frontier. https:/ fpsyg.2	unded ex s in Psystems /www. 2016.0	chology front 0567.	compo ,7,201 Lersin	nents a 16. ISS	pproach SN 1664 /jourr	in the 1078. nals/	
Luyu G Graha <i>Confe</i>	ao, Aman Ma am Neubig. P erence on Ma	adaan, Shuya AL: Program <i>chine Learnii</i>	n Zhou, U -aided lan 1g, pp. 107	ri Alon, F guage mo 764–10799	Pengfei I dels. In 9, 2023.	Liu, Yin Proceed	ning Ya lings of	ng, Jan the 40th	nie Calla h Interno	an, and ational	
Gemma size, 2	Team and G 2024. URL h	oogle DeepN ttps://ar	lind. Gen	n <mark>ma 2: Im</mark> g/abs/2	proving	open la 0118.	inguage	e model	s at a pr	actical	
Edward and V	J. Hu, Yelong Weizhu Chen.	g Shen, Philli LoRA: Low-	p Wallis, Z Rank Ada	Zeyuan Al aptation of	llen-Zhu f Large I	, Yuanzl Languag	hi Li, Si e Mode	hean W els, 202	ang, Lu 1.	Wang,	
Albert (oth'e theory 2023.	Qiaochu Jiang e Lacroix, Yu em provers w . URL http:	, Sean Welled huai Wu, and ith informal p s://doi.o	ck, Jin Per d Guillaur proofs. In rg/10.4	ng Zhou, V ne Lample Internatio 8550/a:	Wenda L e. Draft <i>nal Conj</i> rXiv.2	i, Jiache , Sketch <i>ference</i> 210.1	eng Liu , and P <i>on Lear</i> 2283.	, Mateja Prove: C rning Ra	a Jamnik Juiding Spresent	k, Tim- formal <i>ations</i> ,	
Subbara Acad URL	ao Kambhamj <i>emy of Scienc</i> http://dx	oati. Can lar es, 1534(1):1 doi.org,	ge langua 5–18, Ma /10.111	ge models rch 2024. 1/nyas.	s reason ISSN 1 15125	and pla 749-663	un? Ar. 32. doi:	nals of 10.111	<i>the Ne</i> 1/nyas.	w York 15125.	
Yujia Li Eccle de M Gowa Pushi code 9203 abq1	i, David Choi es, James Kee lasson d'Autu al, Alexey Cl meet Kohli, M generation w . doi: 10.112 1158.	, Junyoung C ling, Felix G une, Igor Ba herepanov, Ja Vando de Fre ith AlphaCod 26/science.ab	Chung, Na imeno, As buschkin, ames Mol itas, Kora e. <i>Science</i> q1158. U	te Kushm gustin Dal Xinyun (loy, Danio y Kavukc g, 378(662 RL http	an, Julia Lago, T Chen, Po el J. Ma uoglu, a 24):1092 o://dx	n Schrif Thomas S-Sen H ankowitz nd Orio –1097, 1 . doi.c	ttwieser Hubert Juang, J z, Esmo l Vinya Deceml Drg/10	r, Rémi r, Peter Johanne e Suthe ils. Con ber 202 0.112	Leblond Choy, C s Welbl rland R npetition 2. ISSN 6/scie	d, Tom Cyprien , Sven obson, n-level 1095- ence.	

Derrick Somerset Macnutt. Ximenes on the art of the crossword. Methuen, 1966.

648 649 650	Tomas Mikolov, Edouard Grave, Piotr Bojanowski, Christian Puhrsch, and Armand Joulin. Advances in pre-training distributed word representations. In <i>Proceedings of the International Conference on Language Resources and Evaluation (LREC 2018)</i> , 2018.
652 653	Tal Ridnik, Dedy Kredo, and Itamar Friedman. Code generation with AlphaCodium: From prompt engineering to flow engineering, 2024.
654 655 656 657	Josh Rozner, Christopher Potts, and Kyle Mahowald. Decrypting cryptic crosswords: Semantically complex wordplay puzzles as a target for NLP. In <i>Advances in Neural Information Processing Systems</i> , volume 34, pp. 11409–11421, 2021. URL https://papers.nips.cc/paper/2021/hash/5fld3986fae10ed2994d14ecd89892d7-Abstract.html.
658 659 660	Abdelrahman "Boda" Sadallah, Daria Kotova, and Ekaterina Kochmar. Are LLMs good cryptic crossword solvers?, 2024. URL https://arxiv.org/abs/2403.12094.
661 662	Soumadeep Saha, Sutanoya Chakraborty, Saptarshi Saha, and Utpal Garain. Language models are crossword solvers, 2024. URL https://arxiv.org/abs/2406.09043.
663 664	unsloth.ai. Unsloth code repo. https://github.com/unslothai/unsloth, 2024.
665 666	Eric Wallace, Nicholas Tomlin, Albert Xu, Kevin Yang, Eshaan Pathak, Matthew Ginsberg, and Dan Klein. Automated crossword solving, 2022.
667 668 669	David Webb. Epic crossword battle expert vs. Times cryptic puzzle #29029. https://youtu. be/N5p4TqdjsHs, 2024.
670 671 672	Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny Zhou. Chain-of-thought prompting elicits reasoning in large language models, 2023. URL https://arxiv.org/abs/2201.11903.
673 674 675 676	Wikipedia. Cryptic crossword — Wikipedia, the free encyclopedia. https://en.wikipedia. org/w/index.php?title=Cryptic_crossword&oldid=1228427465, 2024. [On- line; accessed 1-July-2024].
677 678	Wikipedia contributors. Cryptic crossword - regional variation. https://en.wikipedia. org/wiki/Cryptic_crossword#Regional_variation, 2024.
679 680 681	PW Williams and D Woodhead. Computer assisted analysis of cryptic crosswords. <i>The Computer Journal</i> , 22(1):67–70, 1979.
682	
684	
685	
686	
687	
688	
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700 704	
701	

702 A APPENDIX

A.1 CRYPTIC CROSSWORD BACKGROUND

The following borrows extensively from the description on Wikipedia (2024) (kudos to the authors there), to which we have added wordplay annotations in a notation typical of the FifteenSquare.com website (and in the Wordplay dataset use in this work).

710 A.1.1 BASICS

A cryptic clue leads to its answer only if it is read in the right way. What the clue appears to say when read normally (the surface reading) is usually a distraction with nothing to do with the solution. The challenge is to find the way of reading the clue that leads to the solution.

- 715 A typical clue consists of two parts:
- The straight or definition. This is in essence the same as any non-cryptic crossword clue: a synonym for the answer. It usually exactly matches the part of speech, tense, and number of the answer, and usually appears at the start or end of a clue. For our annotations, the span that encompasses the definition is highlighted using curly braces.
- The cryptic, subsidiary indication or wordplay. This gives the solver some instructions on how to get to the answer in another (less literal) way. The wordplay parts of clues can be obscure, especially to a newcomer, but they tend to utilise standard rules and conventions which become more familiar with practice.
- 724 725

726

727

741

709

Sometimes the two parts of the clue are joined with a link word or phrase such as 'from', 'gives' or 'could be'. One of the tasks of the solver is to find the boundary between the definition and the wordplay, and insert a mental pause there when reading the clue cryptically.

We list below several of the important styles of wordplay that are commonly used, each with an annotated example. For a more comprehensive list, along with an outline of the 'Ximenean principles', please see Wikipedia (2024).

A.1.2 ANAGRAMS

An anagram is a rearrangement of a certain section of the clue to form the answer. This is usually indicated by a codeword which indicates change, movement, breakage or something otherwise amiss.
For example:

```
737 clue: Chaperone shredded corset (6)
738 definition: {Chaperone} shredded corset
739 answer: ESCORT
740 wordplay: (corset)* (*shredded)
```

742 A.1.3 CHARADE

In a charade, the answer is formed by joining individually clued words to make a larger word (namely, the answer). For example:

```
746 clue: Outlaw leader managing money (7)
747 definition: Outlaw leader {managing money}
748 answer: BANKING
749 wordplay: BAN (outlaw) + KING (leader)
```

750 751 A.1.4 CONTAINERS

A container or insertion clue puts one set of letters inside another. For example (also starting to add a little more indirection):

755 clue: Utter nothing when there's wickedness about (5) definition: {utter} nothing when there's wickedness about

```
756
       answer:
                     VOICE
757
                     O (nothing) with VICE (wickedness) around it (about)
       wordplay:
758
759
       A.1.5 DELETIONS
760
       Deletion is a wordplay mechanism which removes some letters of a word to create a shorter word.
761
       For example:
762
763
       clue:
                     Bird is cowardly, about to fly away (5)
764
       definition: {Bird} is cowardly, about to fly away
       answer:
                     RAVEN
765
                     [c]RAVEN (cowardly) - 'C' (i.e. circa, about) (-fly away)
       wordplay:
766
767
       A.1.6 DOUBLE DEFINITION
768
769
       A clue may, rather than having a definition part and a wordplay part, have two definition parts. For
770
       example:
771
                     Not seeing window covering (5)
772
       clue:
       definition: {Not seeing} {window covering}
773
       answer:
                     BLIND
774
                     Double Definition (DD)
       wordplay:
775
776
       A.1.7 HIDDEN WORDS
777
778
       With hidden word clues, the solution itself is written within the clue – either as part of a longer word
779
       or across more than one word. For example:
780
                     Found ermine, deer hides damaged (10)
       clue:
781
       definition: Found ermine, deer hides {damaged}
782
       answer:
                     UNDERMINED
783
       wordplay:
                     [fo]UND ERMINE D[eer] (hides)
784
785
       A.1.8 HOMOPHONES
786
       Homophones are words that sound the same but have different meanings, such as 'night' and
787
       'knight'. Homophone clues always have an indicator word or phrase that has to do with being
788
       spoken or heard. For example:
789
790
       clue:
                     We hear twins shave (4)
791
       definition: We hear twins {shave}
792
       answer:
                     PARE
                     "pair" (twins, "we hear")
793
       wordplay:
794
       A.1.9 REVERSALS
795
796
       A word that gets turned around to make another is a reversal. For example:
797
798
                     Returned beer fit for a king (5)
       clue:
799
       definition: Returned beer {fit for a king}
       answer:
                     REGAL
800
       wordplay:
                     (LAGER) < (beer, <returned)
801
802
803
804
805
806
807
808
809
```

810	A.2	WORDPLAY DATASET
811		

812 The Wordplay Dataset used in this work is extracted from websites where cryptic crossword en-813 thusiasts post solutions to the puzzles published in major publications. Each completed puzzle is annotated by an solver who provides the community with definition, wordplay and answer 814 fields for each of the approximately 30 clues in that day's grid. 815

- 816 For UK papers, these enthusiast websites include: 817
- 818 • timesforthetimes.co.uk - Times, Times Quick
- 819 • www.fifteensquared.net - Independent, Guardian, Financial Times
- 820 bigdave44.com - Telegraph, Sunday Telegraph 821

822 The following is an example from the Wordplay dataset, formatted in YAML (the workings of this clue are illustrated in Figure 2c): 824

```
825
       title: Financial Times 16,479 by FALCON
       url: https://www.fifteensquared.net/2020/05/18/ \
            financial-times-16479-by-falcon/
827
       author: teacow
828
       clues:
829
       - clue: '{Offer} of support also broadcast'
830
         pattern: '8'
831
         ad: D
         answer: PROPOSAL
832
         wordplay: PROP (support) + (ALSO) * (*broadcast)
833
         . . .
834
```

In the above:

- clue is the original clue, as given to solvers, but with the 'regular crossword' definition portion highlighted with curly braces;
- pattern is the number of characters in the answer;
- ad (across/down) is potentially significant, because some clues include directional hints such as 'before' or 'upwards' which are only meaningful if the orientation of the answer within the grid is known;
- answer is the clue's final answer (not known to the solvers before solving); and
- wordplay is an informally annotated explanation of how the clue words act together to logically build the letters in the answer (the resulting grid letters typically being in upper case) - here the * symbol signifies that ALSO is to be anagrammed due to the anagram indicator (broadcast) in the clue.

849 The Wordplay dataset is publicly available as Andrews (2024). Note that care was taken to ensure 850 that the training/validation/test splits follow those of the Cryptonite dataset (and the test set answers 851 are deliberately scrubbed from the retrieved data by the provided scripts, to reduce the chance that 852 they become training data for an over-eager crawling system).

853 854

855

859

860

823

835

836 837

838

839

840

841

842

843

844

845

846

847

848

A.3 CHOICE OF CRYPTONITE VS ROZNER

856 At the start of this paper's research program, the Cryptonite dataset of Efrat et al. (2021) was chosen as being the focus, over the approximately contemporaneous dataset from Rozner et al. (2021) (denoted Rozner here), for the following reasons: 858

- Cryptonite was larger (523k clues, compared to 142k in Rozner)
- Cryptonite consists of clues from The Times and The Telegraph (whereas Rozner is the UK's 861 Guardian). While these are all fine newspapers, it is clear that in the cryptic crossword community 862 (found online via websites for wordplay discussions, or YouTube channels) that The Times is 863 considered the Gold Standard of cryptic crosswords.

004	
864	- Indeed, Connor (2024) - one of the Guardian's own cryptic blog posts - directly states: "The
865	Times hosts an annual crossword-solving competition and it remains, until such time as the
866	Guardian has its own version, the gold standard."
867	- In the authors' view. The Times deserves its role as Gold Standard due to (a) adhering to /
868	upholding the Ximenean standard Macnutt (1966) for what is allowed in clues: (b) doing so for
869	decades: and (c) maintaining high consistency of clue difficulty within puzzles (where solvers
870	frequently complain that the Guardian clues can often be rather haphazard)
871	• The Cryptonite dataset was made available for direct download - even though the licensing is
872	(politely) 'fuzzy', it remains a useable research dataset (and seems unlikely to be challenged
873	by The Times, since it is not possible to reconstruct their full puzzles from the clues given as
874	individual line-items, due to deduplication, for example)
875	- The Bozner dataset required researchers to 'scrape their own data' likely because while the
876	data was being retrieved from a public website the data itself could reasonably be assumed to
877	be convergented. This slight inconvenience had a useful impact (please see below)
878	be copyrighted. This slight inconvenience had a discrut impact (please see below)
070	• Unlike the Cryptonite dataset, the Rozner dataset does not include Across/Down markers for the
019	clues - which makes some of the clues difficult to resolve (for instance EXAMPLE on the paper's
088	first page can only be read correctly if one sees that it is a Down clue - which converts 'up' into a
881	reversal indicator)
882	• The Cryptonite dataset also includes 'is_quick' annotations that show whether a clue was taken
883	from a 'Quick Cryptic' crossword (these clues are typically easier, which enables a further degree
884	of performance analysis).
885	• The Cryptonite dataset splits were set in stone. Pozner, though had a series of splits (random
886	disjoint and 'init'):
887	disjonit, and mit).
888	- The 'random' split was clearly shown to be a poor way of separating train/test due to close
889	overlaps
890	 The 'disjoint' split is similar in spirit to the Cryptonite methodology
891	- The 'Init' split had the additional twist that common prefixes would only be found in their own
892	splits. This had a catchy intuition, although it's not clear from a cryptic cluing perspective
803	whether this has much genuine basis. While there are some prefixes that are common (eg: EX-
093	is easily clued by referring to divorce, etc), the impact seems overall marginal (particularly
094	given the accuracy rate differences reported)
895	
896	Our paper describes a system trained on Cryptonite clue/answer training data, and also (as a com-
897	ponent) the Wordplay dataset (which abides by the Cryptonite splits too).
898	It would be possible to test our existing (Cryptonite trained) system on the Rozner 'Init' test set.
899	However, while Saha et al. (2024) could have the flexibility to run tests on either dataset (since no
900	training was performed), running our current model on the Rozner 'Init' test set would be clearly
901	mis-aligned vis-a-vis the data split.
902	Due there is also a structural masses against as training the mass 2^{2} suctors on the D $_{2}$ and $(U_{1}^{2})^{2}$ with
903	but there is also a structural reason against re-training the paper's system on the Kozner 'Init' split
904	for (specifically) wordplay. The wordplay dataset generation process was guided by the principle of maintaining the Counterplay calles it would be a disaster if Deeper (Init). Wordplay calles were to be
905	maintaining the Cryptointe spins, it would be a disaster it Kosner init <i>worapiay</i> spiits were to be made public. The reason: It is very likely that the Cryptonite test set has a large intersection with the
906	Bozner 'Init' training set (and conversely). As scenes avident from the baseline improvements shown
907	above Open AL likely traine on the Crystonite training set (as they are valeements shown
908	since (as of November 2024) Sala et al. (2024) appears to have released (or re-released) the 'Init'
909	training set under an MIT license, a commercial vendor such as Open AI would be quite within their
010	rights to also train on that. Thus, commercial systems (against which reviewers are forcing academic
011	namers to benchmark) will have been trained on the test sets (without commercial vendors evolicitly
311	'cheating' - they will just be training on all the available training data)
912	encouring and the just be training on an the available training tata).
913	In the authors' judgement, the reasoning paths that are being tested here through the cryptic cross-
914	word task are a prize cultural asset, generated over decades of human effort, and this should not be
915	squandered. Hopefully, this explains the authors' reluctance to dataset-hop : We don't want to make
916	it common to gather and distribute cross-contaminating datasets, specifically Wordplay datasets.
917	

918 A.4 FINE-TUNING PROMPT

920 The following is a verbatim training example used for the fine-tuning of the Gemma2-9B-base 921 model: 922 ### Instruction: 923 Cryptic clue wordplay generation : Given the clue and the answer, \setminus 924 return expert definition and wordplay annotations 925 926 ### Input: clue: "musical and ballet, oddly, that can be avoided" 927 answer: EVITABLE ~ evitable 928 929 ### Response: 930 definition: musical and ballet, oddly, {that can be avoided} 931 wordplay: EVITA (musical) + B[a]L[l]E[t] (ballet, odd letters) 932 933 A.5 IN-CONTEXT LEARNING PROMPTS FOR THE GEMINI LLM 934 The Gemini LLM is prompted in-context with the concatenation of the following sections: 935 936 Cryptic Crossword overview 937 Many-shot wordplay examples 938 939 • Declaration of 'external' Python functions 940 6-shot formalisation demonstration 941 • Actual problem statement (for continuation as a Python proof) 942 • After a verification failure: Error messages for the generated proof, with hints if available, and 943 request to improve iteratively 944 945 The sections of the prompt are described more fully below, note that care was taken to ensure that 946 the chosen terminology was use consistently throughout. 947 948 A.5.1 CRYPTIC CROSSWORD PREAMBLE 949 950 The following is the rubric and wordplay preamble given to the Gemini LLM: 951 A Cryptic crossword question involves using the words in \ 952 the given clue to yield an answer that matches the letter pattern. 953 The clue will provide a definition of the answer, as well \setminus 954 as some 'wordplay' that can also be used to confirm the answer. 955 Expert question solvers write informal 'proofs' using a \ 956 particular format. 957 For the definition, the original clue is annotated with $\$ 958 '{}' to denote where the definition is to be found. 959 For the wordplay, the following conventions are loosely used: 960 * The answer is assembled from the letters in CAPS 961 \star Words in brackets show the origins of letters in CAPS, \setminus often being synonyms, or short forms 962 * Action words are annotated as illustrated: 963 + (ETO N) * (*mad = anagram-signifier) = TONE 964 + (FO OR) < (<back = reversal-signifier) = ROOF 965 + [re]USE (missing = removal-signifier) = USE 966 * DD is a shorthand for 'Double Definition' 967 968 969 970 971

972 A.5.2 MANY-SHOT WORDPLAY EXAMPLES 973

```
974
       Around 20 examples from the Wordplay dataset are included in the in-context prompt:
975
       For example:
976
977
       clue: "arrived with an artist, to get optical device (6)"
978
       definition: arrived with an artist, to get {optical device}
979
       answer: CAMERA
980
       wordplay: CAME (arrived) + RA (artist, short form)
981
       clue: ...
982
983
       A.5.3 EXTERNAL PYTHON DSL FUNCTIONS
984
985
       Domain Specific Python functions are described in-context to the LLM, which appears able to use
986
       them without seeing their internal functionality. In fact, the actual implementation of the functions
       is more extensive than described, since calls to these functions also track 'near misses' which can
987
       be fed back as hints during the re-write process.
988
989
       The task is to produce a formal proof using python code, \setminus
990
       where the docstring will also include an informal proof as an aid.
991
       The following are functions that can be used in your output code:
992
993
      Action=Enum('Action', 'ANAGRAM, REMOVE_FIRST, INITIALS, REMOVE_LAST, '+
                               'GOES_INSIDE, GOES_OUTSIDE, REVERSE, SUBSTRING, HOMOPHONE')
994
       # External definitions
995
       def is_synonym(phrase:str, test_synonym:str, pattern:str='') -> bool:
996
         # Determines whether 'test_synonym' is a reasonable synonym for 'phrase',
997
         # with letters optionally matching 'pattern'
       def is_abbreviation(phrase:str, test_abbreviation:str) -> bool:
998
         # Determines whether 'test abbreviation' is
999
         # a valid abbreviation or short form for 'phrase'
1000
       def action_type(phrase:str, action:Action) -> bool:
1001
         # Determines whether 'phrase' might signify the given 'action'
       def is_anagram(letters:str, word:str) -> bool:
         # Determines whether 'word' can be formed from 'letters' (i.e. an anagram)
1003
       def is_homophone(phrase:str, test_homophone:str) -> bool:
1004
         # Determines whether 'test_homophone' sounds like 'phrase'
1005
1006
       A.5.4 FEW-SHOT FORMALISATION EXAMPLES
1007
1008
       The following are 3 (out of 6) of the few-shot formalisation examples given before the final test-case
1009
       prompt:
1010
       The following are examples of simple functions that prove that \
1011
       each puzzle solution is correct:
1012
1013
       •••• python
1014
       def proof(answer="ONCE",
1015
                  clue="head decapitated long ago", pattern='4'):
         .....
1016
         definition: head decapitated {long ago}
1017
         wordplay: [b]ONCE (head decapitated = remove first letter of BONCE)
1018
1019
         assert is_synonym("head", "BONCE")
1020
         assert action_type("decapitated", Action.REMOVE_FIRST) \
                and "BONCE" [1:] == "ONCE"
1021
         assert is_synonym("long ago", "ONCE", pattern='4')
1022
       proof()
1023
1024
       ··· python
1025
```

```
def proof(answer="DECIMAL",
```

```
1026
                  clue="the point of medical treatment", pattern='7'):
1027
         ....
1028
         definition: {the point} of medical treatment
         wordplay: (MEDICAL) * (*treatment = anagram)
1029
         ......
1030
         assert is_synonym("the point", "DECIMAL", pattern='7')
1031
         assert action_type("treatment", Action.ANAGRAM)
1032
         assert is_anagram("MEDICAL", "DECIMAL")
1033
       proof()
1034
1035
       ··· python
1036
       def proof(answer="SUPERMARKET",
1037
                  clue="fat bags for every brand that's a big seller",
1038
                  pattern='11'):
         ....
1039
         definition: fat bags for every brand that's {a big seller}
1040
         wordplay: SUET (fat) (bags = goes outside) of \
1041
                    (PER (for every) + MARK (brand))
1042
         ....
1043
         assert is_synomym("fat", "SUET")
         assert action_type("bags", Action.IS_OUTSIDE)
1044
         assert "SUET" == "SU" + "ET"
1045
         assert is_abbreviation("for every",
                                                 "PER")
1046
         assert is_synomym("brand", "MARK")
1047
         assert "SU"+"PER"+"MARK"+"ET" == "SUPERMARKET"
1048
         assert is_synonym("a big seller", "SUPERMARKET", pattern='11')
1049
       proof()
1050
1051
       A.5.5 FORMALISATION INSTRUCTION
1052
1053
      The following instruction is given before the final 'test-case' prompt illustrated in Figure 4:
1054
1055
       # Please complete the following in a similar manner, and return the whole function:
1056
       •••• python
1057
      def proof(answer= ...
1058
1059
       A.5.6 PROOF VERIFICATION WITH HINTING
1061
       Examples of assertion failures, with constructive hinting, are shown:
1062
1063
       AssertionError: assert: is_abbreviation('an Artist', 'RA') :
1064
          'an Artist' does not have a valid abbreviation;
1065
          'RA' is an abbreviation for : artist, artillery, Royal Artillery,
1066
          gunners, painter
       AssertionError: assert action_type('goes crazy', Action.ANAGRAM) :
1067
         'goes crazy' itself does not suggest Action.ANAGRAM, but 'crazy' does
1068
       AssertionError: assert action_type('worked', Action.HOMOPHONE) :
1069
         'worked' does not suggest Action.HOMOPHONE, but maybe Action.ANAGRAM
1070
1071
       # Please re-implement the SOLUTION above \
       (altering both the docstring and the python code as required), \setminus
1072
       taking care to fix each of the problems identified, \
       and return the whole function:
1074
1075
       •••• python
       def proof(answer= ...
1076
1077
       Once the prover has fully parsed a given output with zero assertion failures, the proof is considered
1078
       a success (up to 2 re-write iterations are allowed, more that that is considered an overall failure to
1079
       prove the answer).
```

				Validation			Test		
Model	known%	samples	Overall	Quick	Hard	Overall	Quick	Hard	
GPT-4T ('Init')	25%					33.7%			
GPT-4T ('Init')	50%					52.9%			
GPT-4T ('Init')	70%					76.3%			
Gemini-Flash	25%	200	37.0%	38.5%	36.9%	45.5%	66.7%	43.8%	
Gemma2-9B-it	25%	200	37.5%	38.5%	37.4%	44.0%	66.7%	42.2%	
FastText k=1 NN	25%	200	15.5%	15.4%	15.5%	21.0%	33.3%	20.0%	
FastText k=1 NN	50%	200	52.5%	38.5%	53.5%	62.0%	46.7%	63.2%	
FastText k=1 NN	70%	200	79.0%	61.5%	80.2%	81.0%	100.0%	79.5%	

1094

1000

1095

1097 A.5.7 PARTIAL CORRECTNESS METRICS RESULTS

Our results in Table 2, which corresponds to the Exploiting Partially Filled Grids section in Saha et al. (2024), are based on running models on Cryptonite splits, rather than their 'Init'. However, the percentage differences observed are likely large enough outweigh the shift between the two datasets.

The GPT-4T rows in Table 2 are as reported in Saha et al. (2024), and apply to the 'Init' test dataset (i.e. different from our Cryptonite numbers, but still comparable figures in terms of what is being shown here).

The Gemini/Gemma2 rows show the effect of simply filtering the output of our Gemma2-9B finetuned candidate answer proposal model, based on a random letter pattern (using the same formulation as Saha et al. (2024)) and then using the rest of our pipeline. Here, our approach beats the previously reported GPT-4T results, and is itself limited by our first stage Gemma2 model's 'Top-20' candidate answers only containing the correct answer only around 45% of the time.

Our FastText k = 1 kNN systematic approach is clearly very powerful - particularly considering that it does not involve any large model, merely a brute-force search. This only works because the number of known letters for these rows are so high - indeed the 70% level would not be allowed in crosswords that obey the Ximenean guidelines of Macnutt (1966).

If solving complete grids were the target of our research, we would certainly incorporate this kind of solution and overlay the reasoning component to choose from the short-list output (rather than just selecting the first entry, as here). Note that also the performance is bounded above, because the wordlist is not exhaustive - we determined that 7.0% of the gold answers (on the Cryptonite test set) do not appear in the list. This may not be such an issue with the 'Init' dataset, since that wordlist is likely more restricted.

- 1120
- 1121

- 1123
- 1124 1125
- 1126
- 1127
- 1128
- 1129
- 1130
- 1131
- 1132 1133