# DEMONSTRATE–SEARCH–PREDICT:
# Composing retrieval and language models for knowledge-intensive NLP

Omar Khattab [1]     Keshav Santhanam [1]     Xiang Lisa Li [1]     David Hall [1]
Percy Liang [1]     Christopher Potts [1]     Matei Zaharia [1]

## Abstract

Retrieval-augmented in-context learning has emerged as a powerful approach for addressing knowledge-intensive tasks using frozen language models (LM) and retrieval models (RM). Existing work has combined these in simple "retrieve-then-read" pipelines in which the RM retrieves passages that are inserted into the LM prompt. To begin to fully realize the potential of frozen LMs and RMs, we propose DEMONSTRATE–SEARCH–PREDICT (DSP), a framework that relies on passing natural language texts in sophisticated pipelines between an LM and an RM. DSP can express high-level programs that bootstrap pipeline-aware demonstrations, search for relevant passages, and generate grounded predictions, systematically breaking down problems into small transformations that the LM and RM can handle more reliably. We have written novel DSP programs for answering questions in open-domain, multi-hop, and conversational settings, establishing in early evaluations new state-of-the-art in-context learning results and delivering 37–120%, 8–39%, and 80–290% relative gains against the vanilla LM (GPT-3.5), a standard retrieve-then-read pipeline, and a contemporaneous self-ask pipeline, respectively. We release DSP at `https://github.com/stanfordnlp/dsp`.

## 1. Introduction

In-context learning adapts a frozen language model (LM) to tasks by conditioning the LM on a textual prompt including task instructions and a few demonstrating examples (McCann et al., 2018; Radford et al., 2019; Brown et al., 2020). For knowledge-intensive tasks such as question answering, fact checking, and information-seeking dialogue, retrieval models (RM) are increasingly used to augment prompts

[1] **Stanford University**. Correspondence to: **Omar Khattab** <okhattab@cs.stanford.edu>.

*Preprint.*



*Figure 1.* A comparison between three systems based on GPT-3.5 (`text-davinci-002`). On its own, the LM often makes false assertions. An increasingly popular retrieve-then-read pipeline fails when simple search can't find an answer. In contrast, a task-aware DSP program successfully decomposes the problem and produces a correct response. Texts edited for presentation.

with relevant information from a large corpus (Lazaridou et al., 2022; Press et al., 2022; Khot et al., 2022).

Recent work has shown such *retrieval-augmented in-context learning* to be effective in simple "retrieve-then-read" pipelines: a query is fed to the RM and the retrieved passages become part of a prompt that provides context for the LM to use in its response. In this work, we argue that the fact that both LMs and RMs consume (and generate or retrieve) natural language texts creates an opportunity for much more sophisticated interactions between them. Fully realizing this would be transformative: frozen LMs and RMs could serve as infrastructure across tasks, enabling ML- and domain-experts alike to rapidly build grounded AI systems at a high level of abstraction and with lower deployment overheads and annotation costs.

Figure 1 begins to illustrate the power of retrieval-augmented in-context learning, but also the limitations of "retrieve-then-read" (Lazaridou et al., 2022; Izacard et al., 2022). Our query is "How many storeys are in the castle David Gregory inherited?" When prompted to answer this, GPT-3.5 (`text-davinci-002`; Ouyang et al. 2022) makes up a fictitious castle with incorrect attributes, highlighting the common observation that knowledge stored in LM parameters is often unreliable (Shuster et al., 2021; Ishii et al., 2022). Introducing an **RM** component helps, as the **LM** can ground its responses in retrieved passages, but a rigid
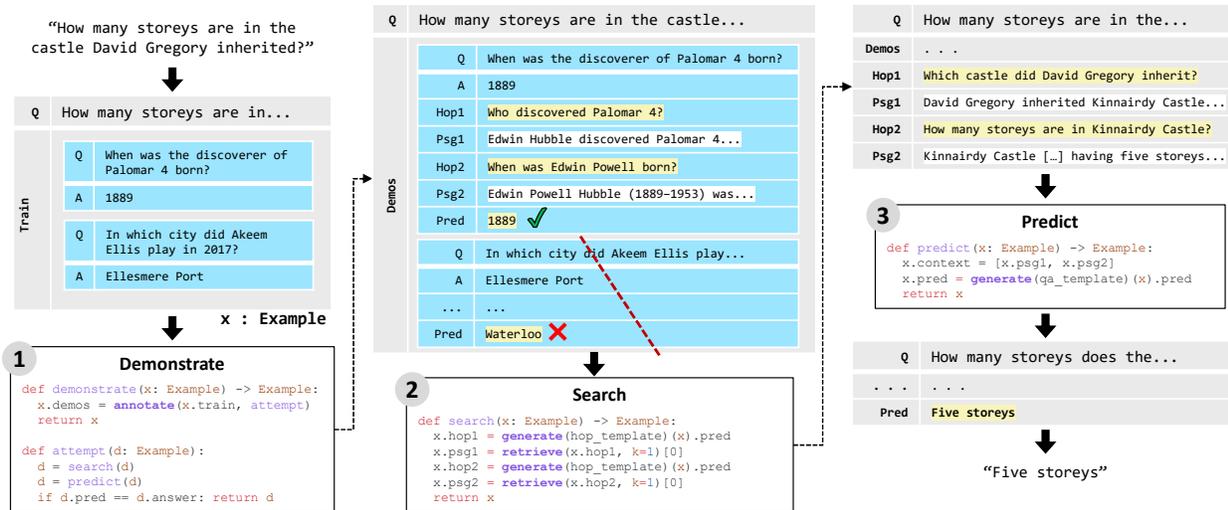
*Figure 2.* A toy example of a DSP program for multi-hop question answering. Given an input question and a 2-shot training set, the DEMONSTRATE stage programmatically annotates intermediate transformations on the training examples using a form of weak supervision. Learning from a resulting *demonstration*, the SEARCH stage decomposes the complex input question and retrieves supporting information over two retrieval hops. Finally, the PREDICT stage uses the demonstration and retrieved passages to answer the question.

retrieve-then-read strategy fails because the **RM** cannot find passages that directly answer the question.

We introduce the **DEMONSTRATE–SEARCH–PREDICT** (**DSP**) framework for in-context learning, which relies entirely on passing natural language text (and scores) between a frozen **RM** and **LM**. DSP introduces a number of composable functions that bootstrap training examples (DEMONSTRATE), gather information from a knowledge corpus (SEARCH), and generate grounded outputs (PREDICT), using them to systematically unify techniques from the retrieval-augmented NLP and the in-context learning literatures (Lee et al., 2019; Khattab et al., 2021a; Anantha et al., 2020; Gao et al., 2022; Izacard et al., 2022; Dohan et al., 2022; Zelikman et al., 2022; Zhang et al., 2022). We use DSP to suggest powerful strategies for knowledge-intensive tasks with compositions of these techniques. This reveals new conceptual possibilities for in-context learning in general (§2), and it allows us to present rich programs that set new state-of-the-art results (§3).

Figure 1 shows the path that a DSP program might take to arrive at an answer, and Figure 2 illustrates how a deliberate program achieves this. Instead of asking the **LM** to answer this complex question, the program's SEARCH stage uses the **LM** to generate a query "Which castle did David Gregory inherit?" The **RM** retrieves a passage saying Gregory inherited the Kinnairdy Castle. After a second search "hop" finds the castle's number of storeys, the PREDICT stage queries the **LM** with these passages to answer the original question. Although this program implements behaviors such as query generation, it requires no hand-labeled examples of these intermediate *transformations* (i.e., of the queries and passages of both retrieval hops). Instead, the DEMONSTRATE

stage uses labeled question–answer pairs to implement a form of weak supervision that programmatically annotates the transformations invoked within SEARCH and PREDICT.

We evaluate several DSP programs on answering questions in open-domain, multi-hop, and conversational settings. In them, we implement novel and reusable transformations such as bootstrapping annotations for all of our pipelines with weak supervision (§2.3), reliably rewriting questions to resolve conversational dependencies and iteratively decompose complex queries with summarization of intermediate hops (§2.4), and generating grounded responses from multiple passages with self-consistency (§2.5). We report preliminary results on Open-SQuAD, HotPotQA, and QReCC using the frozen **LM** GPT-3.5 and **RM** ColBERTv2 (Khattab & Zaharia, 2020; Santhanam et al., 2022b) with no fine-tuning. Our DSP programs deliver 37–120%, 8–39%, and 80–290% relative gains against corresponding vanilla LMs, a standard retrieve-then-read pipeline, and a contemporaneous self-ask pipeline (Press et al., 2022), respectively. Future versions of this report will include additional test tasks and **LM** choices.

In summary, this work makes the following contributions. First, we argue that simple task-agnostic pipelines for in-context learning should give way to deliberate, task-aware strategies. Second, we show that this shift need not be a burden: with DSP, such strategies can be easily expressed as short programs using composable operators. Third, this composability spawns powerful capacities, like automatically annotating demonstrations for complex pipelines from end-task labels. Fourth, for three knowledge-intensive tasks, we implement rich programs that establish state-of-the-art results for in-context learning.

## 2. DEMONSTRATE–SEARCH–PREDICT

We now introduce the DSP framework and show its expressive power by suggesting a number of strategies in which the **LM** and **RM** can come together to tackle complex problems effectively. We show in §3 that such strategies outperform existing in-context learning methods. We begin by discussing the **LM** and **RM** foundation modules on which DSP is built (§2.1) and then the datatypes and control flow within DSP (§2.2). Subsequently, we discuss each of the three inference stages: DEMONSTRATE (§2.3), SEARCH (§2.4), and PREDICT (§2.5).

### 2.1. Pretrained Modules: LM and RM

A DSP program defines the communication between the language model **LM** and the retrieval model **RM**.

**Language Model**  We invoke a frozen language model **LM** to conditionally generate (or score) text. For each invocation, the program prepares a *prompt* that adapts the **LM** to a specific function (e.g., answering questions or generating queries). A prompt often includes instructions, a few demonstrations of the desired behavior, and an input query to be answered.

As in Figure 2, the **LM** generates not only: **(i)** the final answer to the input question (in the PREDICT stage), but also **(ii)** intermediate "hop" queries to find useful information for the input question (SEARCH) as well as **(iii)** exemplar queries that illustrate how to produce queries for questions in the training set (DEMONSTRATE). This systematic use of the **LM** is a hallmark of DSP programs.

**Retrieval Model**  DSP programs also invoke a frozen retrieval model **RM** to retrieve the top-$k$ most "relevant" text sequences for a given *query*. The **RM** can index a massive set of pre-defined *passages* for scalable search, and those passages can be updated without changing the retrieval parameters. The **RM** accepts free-form textual inputs and specializes in estimating the relevance (or similarity) of a text sequence to a query.

As in Figure 2, the **RM** is responsible for retrieving **(i)** passages for each query generated by the **LM** (during the SEARCH stage), but also **(ii)** passages that are used within demonstrations (DEMONSTRATE). In the latter case, the **RM**'s contributions are less about providing directly relevant information to the input question and more about helping the **LM** adapt to the domain and task.

Though not utilized in this example, the **RM** is also used in DSP for functions like retrieving "nearest-neighbor" demonstrations from task training data (DEMONSTRATE) and selecting well-grounded generated sequences from the **LM** (PREDICT).

### 2.2. Datatypes and Control Flow

We have implemented the DSP framework in Python. The present section introduces the core data types and composable functions provided by the framework. We use illustrative code snippets to ground the examples, and to convey the power that comes from being able to express complex interactions between the **LM** and **RM** in simple programs.

**The `Example` Datatype**  To conduct a task, a DSP program manipulates one or more instances of the `Example` datatype. An `Example` behaves like a Python dictionary with multiple fields. The program is typically provided with a few training examples. The code snippet below illustrates this for multi-hop question answering.

```python
from dsp import Example

train = [Example(question="When was the discoverer
    of Palomar 4 born?", answer="1889"),
        Example(question="In which city did Akeem
    Ellis play in 2017?", answer="Ellesmere Port")]
```

This snippet contains two labeled examples, each with a multi-hop question (e.g., "In which city did Akeem Ellis play in 2017?") and its short answer ("Ellesmere Port"). Arbitrary keys and values are allowed within an `Example`, though typical values are strings or lists of strings.

In this task, we are unlikely to find an individual passage that provides the answer to any question. For example, the first training example can probably be resolved only by first answering the question of who discovered Palomar ("Edwin Hubble") and then addressing the question of Hubble's birth date using different evidence passages. We typically assume that the human-labeled training data do *not* include labels for intermediate transformations (e.g., queries for individual hops) that would be useful for following these steps, and so it is the job of the DSP program to discover these strategies via in-context learning.

**A DSP Program**  The following code snippet is a complete program for resolving multi-hop questions like those in Figure 1, with help from train examples like those above.

```python
def multihop_program(question: str) -> str:
    x = Example(question=question, train=train)
    x = multihop_demonstrate(x)
    x = multihop_search(x)
    x = multihop_predict(x)
    return x.answer

multihop_program("How many storeys does the castle
    David Gregory inherited have?")
# => "five storeys"
```

The program takes the input (here, a question) and outputs the system output (its short answer). It starts by creating an `Example` for the input question and assigning the `train` field to the training set from the previous snippet. Programs

invoke and compose DSP *primitives* (i.e., built-in functions) to build the DEMONSTRATE, SEARCH, and PREDICT transformations that define the program.

**Transformations** A transformation is a function that takes an `Example` as input and returns an `Example`, populating new fields (or modifying existing fields) in it. This program invokes three developer-defined transformations, namely, `multihop_demonstrate`, `multihop_search`, and `multihop_predict`. Transformations may themselves invoke other transformations, and they act analogously to layers in standard deep neural network (DNN) programming frameworks such as PyTorch, except that they pass text data instead of tensors between each other and do not involve backpropagation.

We categorize transformations according to their behavior (or purpose) under one of the DEMONSTRATE, SEARCH, and PREDICT stages. That said, DSP does not impose this categorization and allows us to define functions that may blend these stages. We will discuss each of the three stages next.

## 2.3. DEMONSTRATE

It is known that including examples of the desired behavior from the **LM** in its prompt typically leads to better performance (Brown et al., 2020). In DSP, a *demonstration* is a training example that has been prepared to illustrate specific desired behaviors from the **LM**. A DEMONSTRATE transformation takes as input x of type `Example` and prepares a list of demonstrations in `x.demos`, typically by *selecting* a subset of the training examples in `x.train` and *bootstrapping* new fields in them.

**Bootstrapping Demonstrations** Examples in the training set typically consist of the input text and the target output of the task. The DEMONSTRATE stage can augment a training example by programmatically bootstrapping annotations for intermediate transformations. In our running "multi-hop" example, the demonstrations illustrate three **LM**-based transformations: **(i)** how to break down the input question in order to gather information for answering it (i.e., first-hop retrieval), **(ii)** how to use information gathered in an earlier "hop" to ask follow-up questions, and **(iii)** how to use the information gathered to answer complex questions.

```
1 Examples = list[Example]
2 Transformation = Callable[[Example],
3                           Optional[Example]]
4
5 annotate(train: Examples, fn: Transformation)
6     -> Examples
```

Akin to a specialized `map`, the `annotate` primitive accepts a user-defined transformation `fn` and applies it over a list of training examples. Whenever `fn` returns an example (rather than None), `annotate` caches the intermediate predictions (i.e., the generated queries and retrieved passages). These predictions serve as successful demonstrations for the pipeline transformations. In simple uses, `fn` may attempt to answer the example "zero-shot" one or more times. This is typically done by invoking the SEARCH and PREDICT stages of the program. When an answer is produced, if `fn` assesses it as correct, it returns a populated example in which the intermediate predictions are present.

**Case Study** The snippet below defines the function `multihop_demonstrate`, called in Line 3 of `multihop_program`, and illustrates the usage of `annotate`.

```
1 from dsp import sample, annotate
2
3 def attempt_example(d: Example):
4     d = d.copy(demos=[])
5     d = multihop_search(d)
6     d = multihop_predict(d)
7     return d if d.pred == d.answer else None
8
9 def multihop_demonstrate(x: Example):
10    demos = annotate(x.train, attempt_example)
11    return Example(x, demos=demos)
```

In Line 10, `multihop_demonstrate` invokes `annotate`, which bootstraps missing fields in training examples by caching annotations from `attempt_example`. The transformation `attempt_example` takes a training example d and attempts to answer it in a zero-shot fashion: it creates a copy of d with no demonstrations (Line 4; i.e., zero-shot) and invokes the multi-hop search and predict pipeline (Lines 5 and 6). Each transformation returns an updated version of d with additional fields populated. If the pipeline answers correctly (Line 7), the updated d is returned.

Figure 2 illustrates this behavior. DEMONSTRATE transforms a training question–answer pair to a fully-populated demonstration, including fields such as hop1 and hop2 (i.e., queries for multi-hop search) as well as psg1 and psg2. When the **LM** is later invoked to conduct a transformation, say, generating a "second-hop" query during SEARCH, the psg1 field serves as context and the hop2 field serves as a label for this particular training example.

**Discussion** This simple case study illustrates the power of composition in the DSP abstraction. Because the pipeline is a well-defined program in which transformations communicate by passing text attached to `Examples`, a simple map-and-filter strategy can leverage the **LM** and **RM** to bootstrap annotations for a full *pipeline* from end-task labels. This is an extensible strategy, but even in its simplest form it generalizes the approaches explored recently by Zelikman et al. (2022), Wei et al. (2022), Zhang et al. (2022), and Huang et al. (2022) in which an **LM** self-generates chain-of-thought rationales for an individual prompt.

By bootstrapping pipelines, DEMONSTRATE makes it easy to explore complex strategies in SEARCH and PREDICT without writing examples for every transformation. This includes strategies that are challenging to explore without custom annotations in traditional retrieval-augmented NLP. For instance, Khattab et al. (2021a) introduces a pipeline for multi-hop reasoning that is trained with weak supervision, extending work by Lee et al. (2019) and Khattab et al. (2021b). In it, the target 3 or 4 passages that need to be retrieved must be labeled but the system discovers the best *order* of "hops" automatically.

In contrast, DSP allows us to build complex pipelines without labels for intermediate steps, because we can compose programs out of small transformations. If **LM** and **RM** can accurately process such transformations "zero-shot" (i.e., without demonstrations) on at least one or two examples, these examples can be discovered with end-task labels and used as demonstrations.

To draw on our earlier analogy with DNN frameworks like PyTorch, DEMONSTRATE aims to replace the function of backpropagation in extensible ways by simulating the behavior of the program (corresponding to a "forward" pass) and programmatically learning from errors. In doing this with frozen models and with only end-task labels, DEMONSTRATE introduces a high degree of modularity. In particular, without hand-labeling intermediate transformations, developers may swap the training domain, update the training examples, or modify the program's strategy, and use `annotate` to automatically populate all of the intermediate fields for demonstrations.

**Selecting Demonstrations**   It is not always possible to fit all of the training examples in the context window of the **LM**. DSP provides three primitives for selecting a subset of training examples, namely, `sample`, `knn`, and `crossval`.

```
1 sample(train: Examples, k: int)
2     -> Examples
3
4 knn(train: Examples, cast: Callable[[Example], str])
5     -> fn(example: Example, k: int) # currying
6     -> Examples
7
8 crossval(train: Examples, n: int, k: int)
9     -> fn(evaluate: Transformation)
10     -> Examples
```

As a baseline choice, $k$ demonstrations can be randomly sampled from `train` using the `sample` primitive, an approach used by Brown et al. (2020) and much subsequent work. We can also leverage the **RM**'s representations and select from the training set the $k$ nearest neighbors to the input text, a strategy explored by Liu et al. (2021). Another strategy is to apply cross-validation to select among a number of sampled sets of demonstrations (Perez et al., 2021). For example, given $|\texttt{train}| = 100$ training examples, `crossval`

would select $n$ subsets of $k = 5$ examples each, and return the set with which a transformation `evaluate` performs best on the remaining 95 examples.

**Compositions & Extensions**   By manipulating demonstrations and higher-order transformations, these simple selection and bootstrapping primitives can be combined to conduct larger novel strategies. If the training set is very large (e.g., $|\texttt{train}| = 100,000$), we can conduct knn to find the nearest $k = 16$ examples and only `annotate` these, arriving at a system that learns incrementally in real-time. If the training set is moderately large (e.g., $|\texttt{train}| = 1000$), we can conduct `crossval` and cache the performance of all prompts it evaluates on each training example. At test time, we can use knn to find $k = 50$ similar examples to the test input and select the prompt that performs best on these $k$ examples, producing an adaptive system that is informed by the quality of its pipeline on different types of examples.

## 2.4. SEARCH

The SEARCH stage gathers passages to support transformations conducted by the **LM**. We assume a large knowledge corpus—e.g., a snippet of Web, Wikipedia, or arXiv—that is divided into text *passages*. Providing passages to the **LM** facilitates factual responses, enables updating the knowledge store without retraining, and presents a transparency contract: when in doubt, users can check whether the system has faithfully used a reliable source in making a prediction.

In the simplest scenarios, SEARCH can directly query the **RM**, requesting the top-$k$ passages (from a pre-defined index) that match an input question. This baseline instantiation of SEARCH simulates retrieval in most open-domain question answering systems, which implement a "retrieve-then-read" pipeline, like Lee et al. (2019), Khattab et al. (2021b), Lazaridou et al. (2022), and many others.

```
1 from dsp import retrieve
2
3 def simple_search(x):
4     passages = retrieve(query=x.question, k=2)
5     return passages
```

**SEARCH Strategies**   In many scenarios, the complexity of the task demands more sophisticated SEARCH strategies that empower the **RM** to find relevant passages. Our running example (Figure 2) is one such scenario, in which we suspect examples are likely to require *multi-hop* reasoning in particular. Other settings, for instance, pose conversational challenges, in which the information need expressed by a user can only be resolved by taking into account previous turns in the conversation, or demand more extensive planning (Zhong et al., 2022).

In the retrieval-augmented NLP literature, multi-hop search (Xiong et al., 2020; Khattab et al., 2021a) and con-

versational search (Del Tredici et al., 2021; Raposo et al., 2022) pipelines have received much attention. These systems are typically fine-tuned with many hand-labeled query "rewrites" (Anantha et al., 2020), "decompositions" (Geva et al., 2021; Min et al., 2019), or target hops (Yang et al., 2018; Jiang et al., 2020). Supported with automatic annotations from DEMONSTRATE, the SEARCH stage allows us to simulate many such strategies and many others in terms of passing queries, passages, and demonstrations between the **RM** and **LM**. More importantly, SEARCH facilitates our vision of advanced strategies in which the **LM** and **RM** co-operate to incrementally plan a research path for which the **RM** gathers information and the **LM** identifies next steps.

**Case Study**   Let us build on our running multi-hop example as a case study. We can define `multihop_search_v2` (Line 4 in our core program), a slightly more advanced version of the SEARCH transformation from Figure 2. This transformation simulates the iterative retrieval component of fine-tuned retrieval-augmented systems like IRRR (Qi et al., 2020), which reads a retrieved passage in every hop and generates a search query (or a termination condition to stop hopping), and Baleen (Khattab et al., 2021a), which summarizes the information from many passages in each hop for inclusion in subsequent hops.

```
1  from dsp import generate
2
3  def multihop_search_v2(x, max_hops=3):
4    x.hops = []
5
6    for hop in range(max_hops):
7      summary, query = generate(hop_template)(x)
8      x.hops.append((summary, query))
9
10     if query == 'N/A': break
11
12     passages = retrieve(query, k=5)
13     x.context = [summary] + passages
14
15   return x
```

In `multihop_search_v2`, Line 7 calls the `generate` primitive, which invokes the **LM** to produce a query for each retrieval hop. The **LM** is conditioned on a prompt that is prepared using the `hop_template` template. (We discuss prompt templates and the `generate` primitive in §2.5.) Here, this template may be designed to generate a prompt that has the following format (e.g., for the second hop).

```
1  My task is to write a simple query that gathers
      information for answering a complex question. I
      write N/A if the context contains all
      information required.
2
3  {Task demonstrations from x.demos, if any}
4
5  Context: {x.context}
6  Question: {x.question}
7  Summary: Let's summarize the above context.
      __{summary}__
8  Search Query: __{query}__
```

As shown, the **LM** is instructed to read the context retrieved in earlier hops and a complex question. It is then prompted to write: **(i)** a summary of the supplied context and **(ii)** a search query that gathers information for answering that question. The generated text will be extracted and assigned to the `summary` and `query` variables in (`multihop_search_v2`; Line 7). On Line 10, we terminate the hops if the query is "N/A". Otherwise, Line 12 retrieves $k = 5$ passages using the query and Line 13 assigns the `context` for the subsequent hop (or for PREDICT), setting that to include the `summary` of all previous hops as well as the passages retrieved in the final hop so far.

**Comparison with self-ask**   It may be instructive to contrast this multi-hop DSP program with the recent "self-ask" (Press et al., 2022) prompting technique, which we compare against in §3. Self-ask can be thought of as a simple instantiation of DSP's SEARCH stage. In it, the **LM** asks one or more "follow-up questions", which are intercepted and sent to a search engine. The search engine's answers are concatenated into the prompt and are used to answer the question. This is essentially a simplified simulation of IRRR (Qi et al., 2020).

As a general framework, DSP can express ideas like self-ask and many other, more sophisticated pipelines as we discuss in the present section. More importantly, DSP offers a number of intrinsic advantages that lead to large empirical gains: 80%–290% over self-ask. For instance, DSP programs are deeply modular, which among other things means that DSP programs will annotate and construct their own demonstrations. Thus, they can be developed without labeling any of the intermediate transformations (e.g., the queries generated). In addition, the **LM** prompts constructed by DSP get automatically updated to align with the training data and retrieval corpus provided. In contrast, approaches like self-ask rely on a hand-written prompt with hard-coded examples.

Moreover, DSP assigns the control flow to an explicit program and facilitates design patterns that invoke the **LM** (or **RM**) to conduct small transformations. This allows us to build steps that are dedicated to generating one or more retrieval queries, summarizing multiple passages per hop, and answering questions. These steps are individually simpler than the self-ask prompt, yet our multi-hop DSP program deliberately composes them to build richer pipelines that are thus more reliable. In contrast, self-ask delegates the control flow to the **LM** completions, maintaining state within the prompt itself and intercepting follow-up questions to conduct search. We find that this paradigm leads to a "self-distraction" problem (§3.5) that DSP programs avoid.

**Fusing Retrieval Results**   For improved recall and robustness, we can also *fuse* the retrieval across multiple generated queries. Fusion has a long history in information

retrieval (Fox & Shaw, 1994; Xue & Croft, 2013; Kurland & Culpepper, 2018) and sequentially processing multiple queries was explored recently by Gao et al. (2022) for retroactively attributing text generated by LMs to citations. Inspired by these, we include a `fused_retrieval` primitive to DSP to offer a versatile mechanism for interacting with frozen retrievers. It accepts an optional fusion function that maps multiple retrieval lists into one. By default, DSP uses a variant of CombSUM (Fox & Shaw, 1994), assigning each passage the sum of its probabilities across retrieval lists.

To illustrate, the modification below generates $n = 10$ queries for the transformation `multihop_search_v2`.

```
c = generate(hop_template, n=10)(x)
passages = fused_retrieval(c.queries, k=5)
summary = c.summaries[0] # highest-scoring summary
```

**Compositions & Extensions**  To illustrate a simple composition, we can equip a chatbot with the capacity for conversational multi-hop search by combining a query rewriting step, which produces a query that encompasses all of the relevant conversational context, with the multi-hop transformation, as follows.

```
def conversational_multihop_search(x):
    x.question = generate(conv_rewriting_template)(x)
    return multihop_search_v2(x)
```

Similar approaches can be used for correcting spelling mistakes or implementing pseudo-relevance feedback (Cao et al., 2008; Wang et al., 2022a), in which retrieved passages are used to inform a better search query, though this has not been attempted with pretrained LMs to our knowledge.

## 2.5. PREDICT

The PREDICT stage generates the system output using demonstrations (e.g., in `x.demos`) and passages (e.g., in `x.context`). PREDICT tackles the challenges of reliably solving the downstream task, which integrates much of the work on in-context learning in general. Within DSP, it also has the more specialized function of systematically aggregating information across a large number of demonstrations, passages, and candidate predictions.

**Generating Candidates**  Generally, PREDICT has to produce one or more candidate predictions for the end-task. To this end, the basic primitive in PREDICT is `generate`, which accepts a `Template` and (via currying) an `Example` and queries the **LM** to produce one or more completions, as explored earlier in §2.4. A corresponding primitive that uses the **RM** in this stage is `rank`, which accepts a query and one or more passages and returns their relevance scores.

```
Template   # template: an object that can produce
    prompts and parse completions

generate(template: Template)
    -> fn(example: Example)
    -> Completions   # object with keys to access
    extracted preds and scores

rank(query: str, passages: List[str])
    -> List[float]   # object with keys to access
    passage texts and scores
```

A `Template` is an object that can produce prompts, that is, map an `Example` to a string, and extract fields out of completions. For instance, we can map an example `x` that has a question and retrieved passages to the following prompt:

```
My task is to answer questions using Web documents.

{Task demonstrations from x.demos, if any}

Context: {x.passage}
Question: {x.question}
Rationale: Let's think step by step. __{rationale}__
Answer: __{answer}__
```

As this illustrates, the **LM** will be asked to generate a chain-of-thought rationale (CoT; Wei et al. 2022; Kojima et al. 2022) and an answer, and the generated text will be extracted back into the `rationale` and `answer` keys of each completion.

Each invocation to the **LM** can sample multiple candidate predictions. Selecting a "best" prediction is the subject of much work on decoding (Wiher et al., 2022; Li et al., 2022), but a frozen and general-purpose **LM** may not support custom modifications to decoding. Within these constraints, we present several high-level strategies for selecting predictions and aggregating information in DSP via the **LM** and **RM**.

**Selecting Predictions**  Among multiple candidates, we can simply extract the most popular prediction. When a CoT is used to arrive at the answer, this is the self-consistency method of Wang et al. (2022c), which seeks to identify predictions at which multiple distinct rationales arrive.

```
from dsp import generate, majority

def multihop_predict(x):
    candidates = generate(template_qa)(x)
    return x.copy(answer=majority(candidates).answer)
```

DSP generalizes this in two ways. First, we can sample multiple "pipelines of transformations" (PoT) within the program, rather than locally with "chains of thought" (CoT) in one transformation. These chains may even invoke different paths in the program, as illustrated below.

```
1 from dsp import branch
2
3 def pipeline(x):
4   return multihop_predict(multihop_search_v2(x))
5
6 def PoT_program(question: str) -> str:
7   x = Example(question=question, train=train)
8   x = multihop_demonstrate(x)
9
10  candidates = branch(pipeline, n=5, t=0.7)(x)
11  return x.copy(answer=majority(candidates).answer)
```

In the snippet above, Line 10 invokes the primitive `branch` which samples $n$ different PoTs with a high temperature (e.g., $t = 0.7$) and accumulates their intermediate and final predictions. In this example, our pipeline invokes `multihop_search_v2` (§2.4), which applies a variable number of retrieval hops depending on the questions generated, before doing PREDICT. That is, `PoT_program` potentially invokes multiple distinct paths in the program (i.e., with different multi-hop queries and number of hops in each) across branches. It then selects the `majority` answer overall.

DSP generalizes self-consistency in a second way. When sampling our CoTs or PoTs provides multiple candidates, we can select the top-$k$ (e.g., top-4) predictions and then *compare* them directly. For instance, we may prompt the **LM** to compare these choices as MCQ candidates, a transformation for which DEMONSTRATE can automatically prepare exemplars. This effectively simulates the LM recursion of Levine et al. (2022), though unlike their approach it does not require a large training set or updating any (prompt-tuning) weights. One such implementation is illustrated in `openqa_predict` below.

```
1 def openqa_predict(x):
2   preds = generate(template_qa, n=20)(x).answers
3   x.choices = most_common(preds, k=4)
4
5   queries = [f"{x.question} {c}"
6             for c in x.choices]
7
8   x.passages = fused_retrieval(queries)
9   x.answer = generate(TemplateMCQ)(x).answer
10  return x
```

As an alternative comparison approach, we can invoke the **RM** via rank to find the prediction that is most grounded in a retrieved contexts (i.e., most similar to the concatenation of the retrieved passages) or, given an **RM** that can score completions (Krishna et al., 2022), simply the prediction that has the highest score given the prompt.

**Aggregating Information** When only a few demonstrations or passages are selected, we can simply concatenate them all into the prompt. For instance, GPT-3.5 `text-davinci-002` has a context window of 4097 tokens, which we find to be reasonably large for accommodating several (e.g., 3–5) demonstrations, which individually include their own passages and rationales.

To deal with a larger number of demonstrations or passages, we can `branch` in parallel to process individual subsets of the passages or demonstrations and then aggregate the individual answers using one of the scoring methods presented earlier. Indeed, Lewis et al. (2020) and Lazaridou et al. (2022) have explored marginalization as a way to combine scores across passages and Le et al. (2022) ensemble prompts across demonstrations, which can be expressed in this way.

An alternative aggregation strategy is to accumulate information across passages sequentially, rather than independently. This is effectively how our multi-hop approach works (§2.4). Such a strategy has also been employed recently by Gao et al. (2022) for retroactively attributing text generated by LMs to citations. They generate many queries but instead of fusion (§2.4), they run their pipeline on each query and use its outputs to alter the input to subsequent queries.[1]

## 3. Evaluation

We now consider how to implement DSP programs for three diverse knowledge-intensive NLP tasks: open-domain question answering (QA), multi-hop QA, and conversational QA. All of these tasks are "open-domain", in the sense that systems are given a short question or participate in a multi-turn conversation without being granted access to context that answers these questions.

We build and evaluate intuitive compositions of the functions explored in §2 for each task. We show that, despite low development effort, the resulting DSP programs exhibit strong quality and deliver considerable empirical gains over vanilla in-context learning and a standard retrieve-then-read pipeline with in-context learning.

### 3.1. Evaluation Methodology

In this report, we consider one *development dataset* for each of the tasks we consider, namely, the open-domain version of SQuAD (Rajpurkar et al., 2016; Lee et al., 2019), the multi-hop HotPotQA (Yang et al., 2018) dataset in the open-domain "fullwiki" setting, and the conversational question answering QReCC (Anantha et al., 2020; Vakulenko et al., 2022) dataset, which we used for developing the DSP abstractions. We report the validation set accuracy on all three datasets and discuss them in detail §3.5.

Unless otherwise stated, systems are given access to 16-shot training examples, that is, each DSP program can use (up to) 16 questions—or conversations, where applicable—randomly sampled from the respective training set. We

---

[1] Though most of the functionality in this section is implemented, the primitives `branch`, `knn`, and `crossval` are currently work-in-progress.

subsample the validation and test sets to 1000 questions (or 400 conversations, where applicable) and report average quality across five seeds where each seed fixes a single $k$-shot training set of examples. To control the language model API spending budget, each seed processes one fifth of the evaluation examples (e.g., 200 questions per seed, for a total of 1000 unique questions).

We also dedicate held-out *test datasets* (e.g., Open-NaturalQuestions; Kwiatkowski et al. 2019) and *test tasks* (e.g., claim verification, as in FEVER; Thorne et al. 2018) that we only use for evaluating pre-defined DSP programs rather than development. We will include these results in a future version of this report.

### 3.2. Pretrained Modules

**RM**   We use ColBERTv2 (Santhanam et al., 2022b), a state-of-the-art retriever based on late interaction (Khattab & Zaharia, 2020). We choose ColBERTv2 for its highly effective zero-shot search quality and efficient search (Santhanam et al., 2022a). However, our DSP programs are agnostic to how the retriever represents examples or scores passages, so essentially any retriever can be used.

In addition, by making retrieval a first-class construct, DSP allows us to change or update the search index over time. We simulate this in our experiments by aligning each of our datasets with the nearest Wikipedia corpus among the Dec 2016 Wikipedia dump from Chen et al. 2017, the Nov 2017 Wikipedia "abstracts" dump from Yang et al. 2018, and the Dec 2018 Wikipedia dump from Karpukhin et al. 2020.

**LM**   We use the GPT-3.5 (`text-davinci-002`; Brown et al. 2020; Ouyang et al. 2022) language model. Unless otherwise stated, we use greedy decoding when generating $n = 1$ prediction. We sample with temperature $t = 0.7$ when $n > 1$, like related work (Wang et al., 2022c).

### 3.3. Baselines

**Vanilla LM**   The vanilla LM baselines represent the few-shot in-context learning paradigm used by Brown et al. (2020). The open-domain QA and multi-hop QA baselines randomly sample 16 demonstrations (i.e., all of the examples available to each program in our evaluation) from the training set and do not augment these demonstrations with evidence. Similarly, the conversational QA baseline samples four conversations. The vanilla baselines do not search for passages relevant to the input query.

```
1 def vanilla_LM_QA(question: str) -> str:
2     demos = sample(train, k=16)
3     x = Example(question=question, demos=demos)
4     return generate(qa_template)(x).pred
```

**Retrieve-then-Read**   The "retrieve-then-read" baselines use the **RM** to support each example with a potentially relevant passage before submitting the prompt to the **LM**. This emulates the pipelines used by state-of-the-art open-domain question answering systems (Khattab et al., 2021b; Izacard & Grave, 2020; Hofstätter et al., 2022). In conversational QA, we concatenate the first turn and the final question, an approach that we found to perform much better than simply using the final turn. For multi-hop QA, we retrieve and concatenate two passages per question.

```
1 def retrieve_then_read_QA(question: str) -> str:
2     demos = sample(train, k=16)
3     passages = retrieve(question, k=1)
4     x = Example(question=question,
5                 passages=passages,
6                 demos=demos)
7     return generate(qa_template)(x).pred
```

**Self-ask**   We also compare against self-ask (Press et al., 2022), a contemporaneous pipeline that can be thought of as a specific instantiation of DSP's SEARCH stage followed by a simple PREDICT step. For direct comparison with our methods, we modify the self-ask control flow to query the same ColBERTv2 index used in our DSP experiments instead of Google Search. We evaluate two configurations of self-ask. The first uses the original self-ask prompt template, which contains four hand-written demonstrations. In the second configuration, we modify the prompt template to apply a number of changes that we find are empirically useful for HotPotQA.[2]

### 3.4. Proposed DSP Programs

We build on transformations presented in §2. Our programs for all three tasks have the following structure, illustrated for open-domain QA.

```
1 def openqa_program(question: str) -> str:
2     x = Example(question=question, train=train)
3     x = openqa_demonstrate(x)
4     x = openqa_search(x)
5     x = openqa_predict(x)
6     return x.answer
```

The exception is that the conversational QA program,

---

[2]In particular: **(i)** use ColBERTv2-style passages in the hand-crafted demonstrations of self-ask (i.e., instead of the original Google-style snippets), **(ii)** concatenate 16-shot training examples from the task (i.e., question–answer pairs) as a prefix of the prompt, **(iii)** ask the model to generate a short intermediate answer per retrieval step, and **(iv)** explicitly ask the model to generate a follow-up "search query" at each step. We found the final item to be important because self-ask's default prompt often produces follow-up questions that are not self-contained (e.g., "what is the name of the national park?", which is not an informative search query). We also fix the casing in the prompt to be consistent.

*Table 1.* Development results comparing a task-aware DSP program against baseline vanilla LM and retrieve-then-read LM as well as recent and contemporaneous in-context learning approaches with and without retrieval. All of our runs use GPT-3.5 and our retrieval-based rows use ColBERTv2. The results marked with ¶ are collected from related work as of mid-December 2022, and attributed to their individual sources in the main text. As we discuss in the main text, the marked results are not generally apples-to-apples comparisons, since they span a variety of evaluation settings. Nonetheless, we report them here as qualitative reference points.

| | Open-SQuAD | | HotPotQA | | QReCC | |
|---|---|---|---|---|---|---|
| | EM | F1 | EM | F1 | F1 | nF1 |
| **Vanilla LM** | 16.2 | 25.6 | 28.3 | 36.4 | 29.8 | 18.4 |
| **No-retrieval LM SoTA** | 20.2¶ | – | 33.8¶ | 44.6¶ | – | – |
| **Retrieve-then-Read** | 33.8 | 46.1 | 36.9 | 46.1 | 31.6 | 22.2 |
| **Self-ask** w/ ColBERTv2 Search | 9.3 | 17.2 | 25.2 | 33.2 | – | – |
| **+ Refined Prompt** | 9.0 | 15.7 | 28.6 | 37.3 | – | – |
| **Retrieval-augmented LM SoTA** | 34.0¶ | – | 35.1¶ | – | – | – |
| **Task-aware DSP Program** | **36.6** | **49.0** | **51.4** | **62.9** | **35.0** | **25.3** |

`convqa_program`, accepts `turns` (i.e., a list of strings, representing the conversational history) instead of a single `question`. Unless otherwise stated, our programs default to greedy decoding during the DEMONSTRATE stage.

For SEARCH, our open-domain QA program uses the question directly for retrieving $k = 7$ passages and concatenates these passages into our QA prompt with CoT. For PREDICT, it generates $n = 20$ reasoning chains and uses self-consistency (SC; Wang et al. 2022c) to select its final prediction. For DEMONSTRATE, our open-domain QA program uses the following approach, slightly simplified for presentation. In it, the parameter $k = 3$ passed to `annotate` requests annotating only three demonstrations, which will then be used in the prompts.

```
1 def openqa_demonstrate(x: Example) -> Example:
2     demos = sample(x.train, k=16)
3
4     def openqa_attempt(d: Example) -> Example:
5         d.demos = all_but(demos, d)  # all (raw)
      examples different from d
6
7         d = openqa_search(d, k=2)
8         if not passage_match(d): return None  # skip
      examples where search fails
9
10        d = openqa_predict(d, sc=False)
11        if not answer_match(d): return None  # skip
      examples where predict fails
12
13        return d
14
15    x.demos = annotate(demos, openqa_attempt, k=3)
16    return x
```

Our multi-hop program adopts a very similar approach for DEMONSTRATE and PREDICT. For SEARCH, it uses the approach described in §2.4, with the following adjustments. It uses result fusion across $n = 10$ queries per hop and, among the $n$ predictions, uses the summary corresponding to the largest average log-probability. It uses a fixed number of hops for HotPotQA, i.e., two hops. In each prompt (i.e.,

each hop and QA), it concatenates the summaries of all previous hops (i.e., hop 1 onwards) and a total of $k = 5$ passages divided between the hops (i.e., five passages from the first hop or two passages from the first and three from the second).

For conversational QA, we use a simple PREDICT which generates a response with greedy decoding, conditioned on all of the previous turns of the conversation and five retrieved passages. For SEARCH, our conversational QA pipeline generates $n = 10$ re-written queries (and also uses the simple query as the retrieve-and-read baseline; §3.3) and fuses them as in §2.4. We implement DEMONSTRATE similar to `openqa_demonstrate`, but sample only four examples (i.e., four conversational turns; instead of 16 questions as in open-domain QA) for demonstrating the task for the higher-order transformation `convqa_attempt`, which is passed to `annotate` (not shown for brevity).

```
1 def convqa_attempt(d: Example) -> Example:
2     d.demos = all_but(demos, d)  # all (raw)
      examples that don't intersect with the
      conversation of d
3
4     d = convqa_search(d, k=2)
5     if max(precision(d.answer, p) for p in
      d.passages) < .8: return None  # skip examples
      where search fails
6
7     d = convqa_predict(d, n=20)
8     if max(F1(c.pred, d.answer) for c in
      d.candidates) < .75: return None  # skip
      examples where predict fails out of n=20
      attempts
9
10    return d
```

### 3.5. Development Datasets & Results

**Open-SQuAD**   We conduct the open-domain version of SQuAD over the Wikipedia 2016 corpus from Chen et al. (2017), as processed by Khattab et al. (2021b). We use the

same train/validation/test splits as in Karpukhin et al. (2020) and Khattab et al. (2021b).

Table 1 reports the answer EM and F1. The task-aware DSP program achieves 36.6% EM, outperforming the vanilla LM baseline by 126% EM relative gains. This indicates the importance of grounding the **LM**'s predictions in retrieval, and it shows that state-of-the-art retrievers like ColBERTv2 have the capacity to do so off-the-shelf. The proposed DSP program also achieves relative gains of 8% in EM and 6% in F1 over the retrieve-then-read pipeline, highlighting that non-trivial gains are possible by aggregating information across several retrieved passages as we do with self-consistency.

These in-context learning results are competitive with a number of popular fine-tuned systems. For instance, on the Open-SQuAD test set, DPR achieves 29.8% EM, well below our 16-shot DSP program. On the Open-SQuAD dev set, the powerful Fusion-in-Decoder (Izacard & Grave, 2020) "base" approach achieves approximately 36% (i.e., very similar quality to our system) when invoked with five retrieved passages. Nonetheless, with the default setting of reading 100 passages, their system reaches 48% EM in this evaluation. This may indicate that similar gains are possible for our DSP program if the PREDICT stage is made to aggregate information across many more passages.

For comparison, we also evaluate the self-ask pipeline, which achieves 9.3% EM, suggesting that its fixed pipeline is ineffective outside its default multi-hop setting. Studying a few examples of its errors reveals that it often decomposes questions in tangential ways and answers these questions instead. We refer to this behavior of the **LM** as "self-distraction", and we believe it adds evidence in favor of our design decisions in DSP. To illustrate self-distraction, when self-ask is prompted with "When does The Kidnapping of Edgardo Mortara take place?", it asks "What is The Kidnapping of Edgardo Mortara" and then asks when it was published, a tangential question. Thus, self-ask answers "1997", instead of the time The Kidnapping of Edgardo Mortara takes place (1858).

For reference, Table 1 also reports (as No-retrieval LM SoTA) the concurrent in-context learning results from Si et al. (2022) using `code-davinci-002`, who achieve 20.2% EM without retrieval and 34.0% EM with retrieval, albeit on a different sample and split of the SQuAD data. Overall, their approaches are very similar to the baselines we implement (vanilla LM and retrieve-then-read), though their retrieval-augmented approach retrieves (and concatenates into the prompt) 10 passages from a Wikipedia dump.

**HotPotQA**    We use the open-domain "fullwiki" setting of HotPotQA using its official Wikipedia 2017 "abstracts" corpus. The HotPotQA test set is hidden, so we reserve the official validation set for our testing. We sub-divide

the training set into 90%/10% train/validation splits. In the training (and thus validation) split, we keep only examples marked as "hard" in the original dataset, which matches the designation of the official validation and test sets.

We report the final answer EM and F1 in Table 1. On HotPotQA, the task-aware DSP program outperforms the baselines and existing work by very wide margins, exceeding the vanilla LM, the retrieve-then-read baseline, and the self-ask pipeline by 82%, 39%, and 80%, respectively, in EM. This highlights the effectiveness of building up more sophisticated programs that coordinate the **LM** and **RM** for the SEARCH step.

These results may be pegged against the evaluation on HotPotQA in a number of concurrent papers. We first compare with non-retrieval approaches, though our comparisons must be tentative due to variation in evaluation methodologies. Si et al. (2022) achieve 25.2% EM with CoT prompting. With a "recite-and-answer" technique for PaLM-62B (Chowdhery et al., 2022), Sun et al. (2022) achieve 26.5% EM. Wang et al. (2022b) achieve 33.8% EM and 44.6 F1 when applying a self-consistency prompt for PaLM-540B. Next, we compare with a contemporaneous retrieval-based approach: Yao et al. (2022) achieve 35.1% EM using a system capable of searching using a Wikipedia API. All of these approaches trail our task-aware DSP program, which achieves 51.4% EM, by large margins.

**QReCC**    We use QReCC (Anantha et al., 2020) in an open-domain setting over Wikipedia 2018. QReCC does not have an official development set, so we sub-divide the training set into 90%/10% train/validation splits. For the first question in every conversation, we use the rewritten question as the original question often assumes access to a ground-truth document. We also filter low-quality examples from QReCC.[3]

We conduct the QReCC conversations in an auto-regressive manner. At turn $t > 1$ of a particular conversation, the system sees its own responses (i.e., not the ground truth responses) to previous turns of the conversation. We report the novel-F1 metric (nF1; Paranjape et al. 2022), which computes the F1 overlap between the system response and the ground truth while discounting common stopwords and terms present in the question (or earlier questions). The results are shown in Table 1, and follow the same general pattern as SQuAD and HotPotQA.

---

[3]We remove conversations that have one or more empty ground-truth answers and conversations that have only one or two questions. We also find many conversations that include "what other interesting facts are in this article?", which conflict with the open-domain formulation and have no well-defined answer. Hence, we remove any conversation that includes the keywords "other interesting" or "else", which we found to be markers of low quality.

## 4. Conclusion

For a long time, the dominant paradigm for building models in AI has centered around multiplication of tensor representations, and in the deep learning era this has given rise to highly modular (layer-wise) designs that allow for fast development and wide exploration. However, these design paradigms require extensive domain expertise, and even experts face substantial challenges when it comes to combining different pretrained components into larger systems.

The promise of in-context learning is that we can build complex systems from pretrained components using only natural language as the medium for giving systems instructions and, as we argue for, allowing components to communicate with each other. In this new paradigm, the building blocks are pretrained models and the core operations are natural language instructions and operations on natural language texts. If we can realize this potential, then we can broaden participation in AI system development, rapidly prototype systems for new domains, and maximize the value of specialized pretrained components.

In the current paper, we introduced the DEMONSTRATE–SEARCH–PREDICT (DSP) framework for retrieval augmented in-context learning. DSP consists of a number of simple, composable functions for implementing in-context learning systems as deliberate *programs*—instead of end-task prompts—for solving knowledge intensive tasks. We implemented DSP as a Python library and used it to write programs for Open-SQuAD, HotPotQA, and QReCC. These programs deliver substantial gains over previous in-context learning approaches. However, beyond any particular performance number, we argue that the central contribution of DSP is in helping to reveal a very large space of conceptual possibilities for in-context learning in general.

## Acknowledgements

## References

Anantha, R., Vakulenko, S., Tu, Z., Longpre, S., Pulman, S., and Chappidi, S. Open-domain question answering goes conversational via question rewriting. *arXiv preprint arXiv:2010.04898*, 2020.

Brown, T., Mann, B., Ryder, N., Subbiah, M., Kaplan, J. D., Dhariwal, P., Neelakantan, A., Shyam, P., Sastry, G., Askell, A., et al. Language models are few-shot learners. *Advances in neural information processing systems*, 33: 1877–1901, 2020.

Cao, G., Nie, J.-Y., Gao, J., and Robertson, S. Selecting good expansion terms for pseudo-relevance feedback. In *Proceedings of the 31st annual international ACM SIGIR conference on Research and development in information retrieval*, pp. 243–250, 2008.

Chen, D., Fisch, A., Weston, J., and Bordes, A. Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pp. 1870–1879, Vancouver, Canada, 2017. Association for Computational Linguistics. doi: 10.18653/v1/P17-1171. URL https://aclanthology.org/P17-1171.

Chowdhery, A., Narang, S., Devlin, J., Bosma, M., Mishra, G., Roberts, A., Barham, P., Chung, H. W., Sutton, C., Gehrmann, S., et al. Palm: Scaling language modeling with pathways. *arXiv preprint arXiv:2204.02311*, 2022.

Del Tredici, M., Barlacchi, G., Shen, X., Cheng, W., and de Gispert, A. Question rewriting for open-domain conversational qa: Best practices and limitations. In *Proceedings of the 30th ACM International Conference on Information & Knowledge Management*, pp. 2974–2978, 2021.

Dohan, D., Xu, W., Lewkowycz, A., Austin, J., Bieber, D., Lopes, R. G., Wu, Y., Michalewski, H., Saurous, R. A., Sohl-Dickstein, J., et al. Language model cascades. *arXiv preprint arXiv:2207.10342*, 2022.

Fox, E. A. and Shaw, J. A. Combination of multiple searches. *NIST special publication SP*, 243, 1994.

Gao, L., Dai, Z., Pasupat, P., Chen, A., Chaganty, A. T., Fan, Y., Zhao, V. Y., Lao, N., Lee, H., Juan, D.-C., et al. Attributed text generation via post-hoc research and revision. *arXiv preprint arXiv:2210.08726*, 2022.

Geva, M., Khashabi, D., Segal, E., Khot, T., Roth, D., and Berant, J. Did aristotle use a laptop? a question answering

benchmark with implicit reasoning strategies. *Transactions of the Association for Computational Linguistics*, 9: 346–361, 2021.

Hofstätter, S., Chen, J., Raman, K., and Zamani, H. Fidlight: Efficient and effective retrieval-augmented text generation. *arXiv preprint arXiv:2209.14290*, 2022.

Huang, J., Gu, S. S., Hou, L., Wu, Y., Wang, X., Yu, H., and Han, J. Large language models can self-improve. *arXiv preprint arXiv:2210.11610*, 2022.

Ishii, Y., Madotto, A., and Fung, P. Survey of hallucination in natural language generation. *ACM Comput. Surv*, 1(1), 2022.

Izacard, G. and Grave, E. Leveraging passage retrieval with generative models for open domain question answering. *arXiv preprint arXiv:2007.01282*, 2020.

Izacard, G., Lewis, P., Lomeli, M., Hosseini, L., Petroni, F., Schick, T., Dwivedi-Yu, J., Joulin, A., Riedel, S., and Grave, E. Few-shot learning with retrieval augmented language models. *arXiv preprint arXiv:2208.03299*, 2022.

Jiang, Y., Bordia, S., Zhong, Z., Dognin, C., Singh, M., and Bansal, M. HoVer: A dataset for many-hop fact extraction and claim verification. In *Findings of the Association for Computational Linguistics: EMNLP 2020*, pp. 3441–3460, Online, 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020. findings-emnlp.309. URL https://aclanthology.org/2020.findings-emnlp.309.

Karpukhin, V., Oguz, B., Min, S., Lewis, P., Wu, L., Edunov, S., Chen, D., and Yih, W.-t. Dense passage retrieval for open-domain question answering. In *Proceedings of the 2020 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pp. 6769–6781, Online, 2020. Association for Computational Linguistics. doi: 10.18653/v1/2020.emnlp-main.550. URL https://aclanthology.org/2020.emnlp-main.550.

Khattab, O. and Zaharia, M. Colbert: Efficient and effective passage search via contextualized late interaction over BERT. In Huang, J., Chang, Y., Cheng, X., Kamps, J., Murdock, V., Wen, J., and Liu, Y. (eds.), *Proceedings of the 43rd International ACM SIGIR conference on research and development in Information Retrieval, SIGIR 2020, Virtual Event, China, July 25-30, 2020*, pp. 39–48. ACM, 2020. doi: 10.1145/3397271.3401075. URL https://doi.org/10.1145/3397271.3401075.

Khattab, O., Potts, C., and Zaharia, M. Baleen: Robust Multi-Hop Reasoning at Scale via Condensed Retrieval. In *Thirty-Fifth Conference on Neural Information Processing Systems*, 2021a.

Khattab, O., Potts, C., and Zaharia, M. Relevance-guided supervision for openqa with ColBERT. *Transactions of the Association for Computational Linguistics*, 9:929–944, 2021b.

Khot, T., Trivedi, H., Finlayson, M., Fu, Y., Richardson, K., Clark, P., and Sabharwal, A. Decomposed prompting: A modular approach for solving complex tasks. *arXiv preprint arXiv:2210.02406*, 2022.

Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. Large language models are zero-shot reasoners. *arXiv preprint arXiv:2205.11916*, 2022.

Krishna, K., Chang, Y., Wieting, J., and Iyyer, M. Rankgen: Improving text generation with large ranking models. *arXiv preprint arXiv:2205.09726*, 2022.

Kurland, O. and Culpepper, J. S. Fusion in information retrieval: Sigir 2018 half-day tutorial. In *The 41st International ACM SIGIR Conference on Research & Development in Information Retrieval*, pp. 1383–1386, 2018.

Kwiatkowski, T., Palomaki, J., Redfield, O., Collins, M., Parikh, A., Alberti, C., Epstein, D., Polosukhin, I., Devlin, J., Lee, K., Toutanova, K., Jones, L., Kelcey, M., Chang, M.-W., Dai, A. M., Uszkoreit, J., Le, Q., and Petrov, S. Natural questions: A benchmark for question answering research. *Transactions of the Association for Computational Linguistics*, 7:452–466, 2019. doi: 10.1162/tacl_a_00276. URL https://aclanthology.org/Q19-1026.

Lazaridou, A., Gribovskaya, E., Stokowiec, W., and Grigorev, N. Internet-augmented language models through few-shot prompting for open-domain question answering. *arXiv preprint arXiv:2203.05115*, 2022.

Le, N. T., Bai, F., and Ritter, A. Few-shot anaphora resolution in scientific protocols via mixtures of in-context experts. *arXiv preprint arXiv:2210.03690*, 2022.

Lee, K., Chang, M.-W., and Toutanova, K. Latent retrieval for weakly supervised open domain question answering. In *Proceedings of the 57th Annual Meeting of the Association for Computational Linguistics*, pp. 6086–6096, Florence, Italy, 2019. Association for Computational Linguistics. doi: 10.18653/v1/P19-1612. URL https://aclanthology.org/P19-1612.

Levine, Y., Dalmedigos, I., Ram, O., Zeldes, Y., Jannai, D., Muhlgay, D., Osin, Y., Lieber, O., Lenz, B., Shalev-Shwartz, S., et al. Standing on the shoulders of giant frozen language models. *arXiv preprint arXiv:2204.10019*, 2022.

Lewis, P. S. H., Perez, E., Piktus, A., Petroni, F., Karpukhin, V., Goyal, N., Küttler, H., Lewis, M., Yih, W., Rocktäschel, T., Riedel, S., and Kiela, D. Retrieval-Augmented

Generation for Knowledge-Intensive NLP Tasks. In Larochelle, H., Ranzato, M., Hadsell, R., Balcan, M., and Lin, H. (eds.), *Advances in Neural Information Processing Systems 33: Annual Conference on Neural Information Processing Systems 2020, NeurIPS 2020, December 6-12, 2020, virtual*, 2020. URL https://proceedings.neurips.cc/paper/2020/hash/6b493230205f780e1bc26945df7481e5-Abstract.html.

Li, X. L., Holtzman, A., Fried, D., Liang, P., Eisner, J., Hashimoto, T., Zettlemoyer, L., and Lewis, M. Contrastive decoding: Open-ended text generation as optimization. *arXiv preprint arXiv:2210.15097*, 2022.

Liu, J., Shen, D., Zhang, Y., Dolan, B., Carin, L., and Chen, W. What makes good in-context examples for gpt-3? *arXiv preprint arXiv:2101.06804*, 2021.

McCann, B., Keskar, N. S., Xiong, C., and Socher, R. The natural language decathlon: Multitask learning as question answering. arXiv:1806.08730, 2018. URL https://arxiv.org/abs/1806.08730.

Min, S., Zhong, V., Zettlemoyer, L., and Hajishirzi, H. Multi-hop reading comprehension through question decomposition and rescoring. *arXiv preprint arXiv:1906.02916*, 2019.

Ouyang, L., Wu, J., Jiang, X., Almeida, D., Wainwright, C. L., Mishkin, P., Zhang, C., Agarwal, S., Slama, K., Ray, A., et al. Training language models to follow instructions with human feedback. *arXiv preprint arXiv:2203.02155*, 2022.

Paranjape, A., Khattab, O., Potts, C., Zaharia, M., and Manning, C. D. Hindsight: Posterior-guided Training of Retrievers for Improved Open-ended Generation. In *International Conference on Learning Representations*, 2022. URL https://openreview.net/forum?id=Vr_BTpw3wz.

Perez, E., Kiela, D., and Cho, K. True few-shot learning with language models. *Advances in Neural Information Processing Systems*, 34:11054–11070, 2021.

Press, O., Zhang, M., Min, S., Schmidt, L., Smith, N. A., and Lewis, M. Measuring and narrowing the compositionality gap in language models. *arXiv preprint arXiv:2210.03350*, 2022.

Qi, P., Lee, H., Sido, O., Manning, C. D., et al. Retrieve, rerank, read, then iterate: Answering open-domain questions of arbitrary complexity from text. *arXiv preprint arXiv:2010.12527*, 2020. URL https://arxiv.org/abs/2010.12527.

Radford, A., Wu, J., Child, R., Luan, D., Amodei, D., Sutskever, I., et al. Language models are unsupervised multitask learners. *OpenAI blog*, 1(8):9, 2019.

Rajpurkar, P., Zhang, J., Lopyrev, K., and Liang, P. SQuAD: 100,000+ questions for machine comprehension of text. In *Proceedings of the 2016 Conference on Empirical Methods in Natural Language Processing*, pp. 2383–2392, Austin, Texas, 2016. Association for Computational Linguistics. doi: 10.18653/v1/D16-1264. URL https://aclanthology.org/D16-1264.

Raposo, G., Ribeiro, R., Martins, B., and Coheur, L. Question rewriting? assessing its importance for conversational question answering. In *European Conference on Information Retrieval*, pp. 199–206. Springer, 2022.

Santhanam, K., Khattab, O., Potts, C., and Zaharia, M. PLAID: An Efficient Engine for Late Interaction Retrieval. *arXiv preprint arXiv:2205.09707*, 2022a.

Santhanam, K., Khattab, O., Saad-Falcon, J., Potts, C., and Zaharia, M. ColBERTv2: Effective and efficient retrieval via lightweight late interaction. In *Proceedings of the 2022 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies*, pp. 3715–3734, Seattle, United States, July 2022b. Association for Computational Linguistics. doi: 10.18653/v1/2022.naacl-main.272. URL https://aclanthology.org/2022.naacl-main.272.

Shuster, K., Poff, S., Chen, M., Kiela, D., and Weston, J. Retrieval augmentation reduces hallucination in conversation. *arXiv preprint arXiv:2104.07567*, 2021.

Si, C., Gan, Z., Yang, Z., Wang, S., Wang, J., Boyd-Graber, J., and Wang, L. Prompting gpt-3 to be reliable. *arXiv preprint arXiv:2210.09150*, 2022.

Sun, Z., Wang, X., Tay, Y., Yang, Y., and Zhou, D. Recitation-augmented language models. *arXiv preprint arXiv:2210.01296*, 2022.

Thorne, J., Vlachos, A., Christodoulopoulos, C., and Mittal, A. FEVER: a large-scale dataset for fact extraction and VERification. In *Proceedings of the 2018 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, Volume 1 (Long Papers)*, pp. 809–819, New Orleans, Louisiana, 2018. Association for Computational Linguistics. doi: 10.18653/v1/N18-1074. URL https://aclanthology.org/N18-1074.

Vakulenko, S., Kiesel, J., and Fröbe, M. SCAI-QReCC shared task on conversational question answering. In *Proceedings of the Thirteenth Language Resources and Evaluation Conference*, pp. 4913–4922, Marseille, France,

June 2022. European Language Resources Association. URL https://aclanthology.org/2022.lrec-1.525.

Wang, X., Macdonald, C., Tonellotto, N., and Ounis, I. Colbert-prf: Semantic pseudo-relevance feedback for dense passage and document retrieval. *ACM Transactions on the Web*, 2022a.

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., and Zhou, D. Rationale-augmented ensembles in language models. *arXiv preprint arXiv:2207.00747*, 2022b.

Wang, X., Wei, J., Schuurmans, D., Le, Q., Chi, E., and Zhou, D. Self-consistency improves chain of thought reasoning in language models. *arXiv preprint arXiv:2203.11171*, 2022c.

Wei, J., Wang, X., Schuurmans, D., Bosma, M., Chi, E., Le, Q., and Zhou, D. Chain of thought prompting elicits reasoning in large language models. *arXiv preprint arXiv:2201.11903*, 2022.

Wiher, G., Meister, C., and Cotterell, R. On decoding strategies for neural text generators. *arXiv preprint arXiv:2203.15721*, 2022.

Xiong, W., Li, X. L., Iyer, S., Du, J., Lewis, P., Wang, W. Y., Mehdad, Y., Yih, W.-t., Riedel, S., Kiela, D., et al. Answering complex open-domain questions with multi-hop dense retrieval. *arXiv preprint arXiv:2009.12756*, 2020. URL https://arxiv.org/abs/2009.12756.

Xue, X. and Croft, W. B. Modeling reformulation using query distributions. *ACM Transactions on Information Systems (TOIS)*, 31(2):1–34, 2013.

Yang, Z., Qi, P., Zhang, S., Bengio, Y., Cohen, W. W., Salakhutdinov, R., and Manning, C. D. Hotpotqa: A dataset for diverse, explainable multi-hop question answering. *arXiv preprint arXiv:1809.09600*, 2018.

Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K., and Cao, Y. React: Synergizing reasoning and acting in language models. *arXiv preprint arXiv:2210.03629*, 2022.

Zelikman, E., Wu, Y., and Goodman, N. D. Star: Bootstrapping reasoning with reasoning. *arXiv preprint arXiv:2203.14465*, 2022.

Zhang, Z., Zhang, A., Li, M., and Smola, A. Automatic chain of thought prompting in large language models. *arXiv preprint arXiv:2210.03493*, 2022.

Zhong, V., Shi, W., Yih, W.-t., and Zettlemoyer, L. Romqa: A benchmark for robust, multi-evidence, multi-answer question answering. *arXiv preprint arXiv:2210.14353*, 2022.