000 DISCIPLE: LEARNING INTERPRETABLE PROGRAMS 001 FOR SCIENTIFIC DISCOVERY 002 003

Anonymous authors

Paper under double-blind review

ABSTRACT

Creating hypotheses for new observations is a key step in the scientific process of understanding a problem in any domain. A good hypothesis that is interpretable, reliable (good at predicting unseen observations), and data-efficient; is useful for scientists aiming to make novel discoveries. This paper introduces an automatic way of learning such interpretable and reliable hypotheses in a dataefficient manner. We propose DiSciPLE (Discovering Scientific Programs using LLMs and Evolution) an evolutionary algorithm that leverages common sense and prior knowledge of large language models (LLMs) to create hypotheses as Python programs for scientific problems. Additionally, we propose two improvements: a program critic and a program simplifier to further improve our method to produce good hypotheses. We evaluate our method on four different real-world tasks in two scientific domains and show significantly better results. For example, we can learn programs with 35% lower error than the closest non-interpretable baseline for population density estimation.

024 025 026

027

041

043

045

004

010 011

012

013

014

015

016

017

018

019

021

023

INTRODUCTION 1

028 The scientific process involves the collection of data, the creation of hypotheses to explain the data, and experimental testing and validation of the hypothesis to produce scientific laws. This process is 029 manual, laborious and time consuming. With advances in LLMs, there is an opportunity to automate the scientific pipeline to greatly accelerate scientific discovery and progress. While prior work in 031 automating science has looked at efficient data collection (Van Horn et al., 2015; Sener & Savarese, 2018), the possibility of automatic hypothesis generation from observational data has been less 033 explored. This is especially important because in many domains like earth science, there are vast 034 troves of data waiting to be analyzed. A method that could automatically scour through this data to identify hypotheses for scientists to explore can significantly speed up the scientists' workflow. Our goal is to address this gap and produce an automatic hypothesis generation framework (fig. 1). 037

What are good desiderata for an automatic hypotheses generation framework? First, clearly the generated hypothesis should fit the observations well. This requires that the framework should be able to search through a large and *expressive* search space. Second, the generated hypotheses should 040 be *reliable*, in that it should correctly predict data even outside of the observations used to generate the hypotheses. Third, a good hypothesis is one that is *interpretable* by scientists, since only then 042 can scientists understand what is being discovered, and frame experiments to test the hypothesis. Finally, the hypothesis generation framework should be sample *efficient*. In other words, scientists 044 coming up with a hypothesis should not need to collect a lot of observations to come up with a good



Figure 1: Our method asks an expert for a problem of interest and the observation data for that problem. Using this information it comes up with interpretable programs as hypothesis explaining the observation data.

hypothesis. In this paper, we ask: *How can we automatically discover expressive, interpretable and reliable hypotheses in a sample-efficient way?*

Existing approaches that discover hypotheses or formulae are not up to the task. Interpretable ML 057 techniques such as concept-bottlenecks (Koh et al., 2020) can produce simple bag-of-words programs for visual classification but are not expressive enough to produce more sophisticated hypotheses. At the other end of the spectrum, black-box neural networks are expressive but offer no insight 060 to the scientists for exploration. Symbolic regression (Cranmer, 2023) techniques are interpretable 061 and can produce sophisticated formulae, but are very slow to converge because of the large search 062 space. This issue is particularly severe in scientific domains where the hypothesis may reference 063 any concept or variable in the scientific vocabulary, making the search space effectively unbounded. Neurosymbolic program learning (Johnson et al., 2017) address this search space with reinforcement 064 learning, but this requires exceptional amount of training data and is not sample-efficient. 065

066 To be able to search through a large and expressive search space, recent works have looked at lever-067 aging the common sense abilities of large language models (LLMs) (Surís et al., 2023). However, 068 while zero-shot generation from LLMs can work well for well-known problems, they are insufficient 069 for scientific domains with novel problems and data. Recent work has looked at leveraging LLMs as the mutation operator in evolutionary algorithms (Romera-Paredes et al., 2024; Chiquier et al., 071 2024). We leverage the common sense and prior knowledge of LLMs in tandem with an evolutionary algorithm-based symbolic regression method to generate programs as hypotheses, resulting in 072 a significantly smaller search space. The interpretable programs are learned on a set of lower-level 073 primitives, including neural modules such as open-world segmentation. Our framework produces 074 neuro-symbolic programs as good hypotheses. 075

076 While performing crossover and mutation with LLMs can produce meaningful programs, several 077 improvements can be made on the framework to further improve the search. First, unlike symbolic regression, the hypotheses evaluation in the evolution loop does not have to be limited to a single regression score. The LLM can leverage more information about the current hypothesis to better guide 079 the search. Therefore, we propose to add a *critic* to our framework that can provide fine-grained information about the hypothesis to better guide the search. Our critic divides the training observation 081 set into fine-grained categories and can provide feedback on categories that the hypothesis does well on or does badly on. This results in programs that can have an all-round better performance. Second, 083 while the LLM is good at exploring the search space, many programs produced by the framework 084 contain frivolous variables, that may not aid the hypotheses. Having such redundant features/parts 085 in programs can result in more complex final hypotheses as well as a larger search space during evolution. Therefore, we also add an analytical (non-llm-based) program simplification method, that 087 can remove such frivolous parts of the program.

We run experiments in two different expert domains (demography and climate science) on four real-world problems and show that we can produce good hypotheses by leveraging LLMs. In all of these cases, the programs generated are *reliable* and result in better generalization than deep neural nets. Moreover, in a few cases, we also observe better performance on even unseen in-distribution data.

Our contributions are:

- We introduce a novel framework **DiSciPLE** (**Di**scovering **Sci**entific **P**rograms using LLMs and **E**volution), that can produce interpretable, reliable, and sample-efficient hypotheses for diverse scientific applications.
- We present two key components: a critic and a program simplification method to DiSci-PLE that can further improve the evolutionary search resulting in better hypotheses.
- We show application of DiSciPLE on real-world problems by applying it on two scientific domains and four different problems and show that our learned hypotheses are indeed more interpretable, reliable, and data-efficient compared to baselines.
- 102 103

094

095

096

098

099

- 2 RELATED WORKS
- 104 105

Concept bottlenecks. Concept bottleneck (Koh et al., 2020) is an approach used to create interpretable-yet-powerful classifiers. The key idea is to train a deep model to predict a set of low-level concepts or bottlenecks and then learn a linear classifier. Such concept bottlenecks have

108 the basis of methods in several areas such as fine-grained recognition (Ferrari & Zisserman, 2007; 109 Huang et al., 2016; Zhou et al., 2018; Tang et al., 2020) and zero-shot learning (Lampert et al., 2013; 110 Akata et al., 2015; Kodirov et al., 2017). However in order to train these models, expensive data is 111 needed to be collected for the bottleneck concepts themselves. One way to reduce this annotation 112 cost is to sequentially ask question in an information theoretically optimized way (Chattopadhyay et al., 2024a;b). Researchers have also automated this pipeline by using large-language models as a 113 knowledge base to propose concept bottleneck models (Menon & Vondrick, 2023; Pratt et al., 2023; 114 Han et al., 2023). Chiquier et al. (2024) proposed an evolutionary algorithm with LLMs as the muta-115 tion operation to discover interpretable concept bottleneck models without prior information. While 116 these models are interpretable, they are very simple in terms of expressive power. In this work, we 117 instead evolve programs, which are more expressive than bag of words, while being interpretable. 118

119 **Symbolic regression.** Symbolic regression (Cranmer, 2023) is a powerful technique for learning 120 such programs on evolution through evolutionary search. Several methods have been proposed to 121 improve the search (Makke & Chawla, 2024), however most symbolic regression techniques cannot 122 solve problems beyond simple mathematical formula. This is partly because the search space of 123 possible hypotheses is too large due to combinatorial explosion. Like ours, recent work in symbolic 124 regression (Grayeli et al., 2024) have also looked at using LLMs to better guide the search. However, 125 they also only test this method on mathematical formula, with limited set of primitives. In real-world 126 problems, the primitives functions are more complex than mathematical operations and can even be 127 open-world for example a text-to-image segmentor.

128

Neuro-Symbolic Program Learning. (Mao et al., 2019; Dong et al.) is another avenue for learning programs as hypotheses for observation datasets or answering questions. These methods typically try to learn both discrete program structures together with neural networks. However, since this optimization is non-differentiable these methods require reinforcement learning (Johnson et al., 2017) or complex non-differentiable optimization techniques (Ellis et al., 2021). The hard optimization issue makes the problem of learning programs sample inefficient in real-world settings. We alternatively use LLMs ability to program to better guide the search for such programs.

Program synthesis with LLMs. Several works have utilized LLM coding ability in different applications such as VQA (Surís et al., 2023; Gupta & Kembhavi, 2022) and robot manipulation (Liang et al., 2023). While the zero-shot inferred code work very well on domains well-known to the internet, they tend to perform poorly on problems in scientific domains, as shown by our results.

Scientific applications. Researchers in numerous scientific domains have looked machine learning tools to build predictive models/hypotheses for their quantities of interest. In this work, we experiment in two such scientific domain of: demography and climate Science. In demography, we focus on the problems of socioeconomic indicator prediction (Yong & Zhou, 2024), namely population density and poverty estimation (Metzger et al., 2024; Xie, 2017). Similarly in climate science we focus on two problems of aboveground biomass prediction (AGB) (Nathaniel et al., 2023) and Contiguous Solar Induced Chlorophyll Fluorescence(CSIF) forecasting (Zhang et al., 2018).

148 149 150

3 Methodology

Our key contribution is a hypothesis generation framework that leverages LLMs and performs evolutionary search. In section 3.1 we formalize the problem of hypothesis generation. In section 3.2 we present our method of incorporating LLMs in the evolutionary search framework. Finally, In section 3.3, 3.4 and 3.5, we discuss the improvements to this framework to speed-up the search.

155 156

157

3.1 PROBLEM FORMULATION

Data/Observation: To come up with hypotheses, scientists first collect data/observations for a problem they are trying to understand. We represent a collection of n such observations using $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$. Here, $x_i \in X$ is the input data/covariates to the hypothesis, and $y_i \in Y$ is the true observed output for a quantity of interest. For the problem of "population density estimation from satellite images", X could be the set of satellite images for different regions



Figure 2: Overview of our evolutionary algorithm with *critic* and *simplification*. We start with an initialized bank of program trying to solve a task. From this bank we sample pairs of programs based on their fitness score and perform crossover/mutations over them to produce new programs. The generated program is further improved by passing it through a critic and then an analytical simplification step. This program is then evaluated and put in the next generation of program bank. The evaluation score of the program is used to determine the fitness for the next iteration of evolution.

185

186

187 188

189

190

191

195 196 and Y would be population density maps for the corresponding satellite images collected through census. Our goal is to discover an interpretable hypothesis h that can explain observations \mathcal{D} . Since we want our hypothesis to be data-efficient, n is typically small: in the order of a few thousand.

Objective: The scientist or expert also provides a natural language name or description descr of the quantity of interest Y. For example, this may be the phrase "population density". As we will see below, this information will be useful in guiding the LLM to search the space well.

192 **Metric:** The scientists or domain experts also provide a hypothesis evaluation metric \mathcal{M} . A good hypothesis can lead to the best evaluation score

$$s(h; \mathcal{D}) = \frac{1}{n} \sum_{i=1}^{n} \mathcal{M}(h(x_i), y_i)$$
(1)

For example, in the case of population density a good metric used by domain experts is L2 error over log *i.e.* $\mathcal{M}(y', y) = ||log(y') - log(y)||_2$ (Metzger et al., 2022; 2024).

200 **Primitives:** Finally, the experts also provide our framework with a set of primitive variables and functions $\mathcal{F} = \{f_1, f_2 \dots f_k\}$, that can be used to construct a hypothesis h. This set of primitive 201 variables can even be unbounded or open-vocabulary. For example, for scientists analyzing remote-202 sensing image data, these primitives can include an open-vocabulary segmentation model (Mall 203 et al., 2023). The primitives also include mathematical, logical, and image operations that can be 204 composed to create powerful hypotheses. Examples of such primitives could be logarithm operation, 205 elementwise maximum, or a distance transform respectively. Finally, the primitives can also be used 206 to query specific attributes of the data, such as the average temperature or average precipitation at 207 a location. The API of primitive functions can also change depending on the application domain. 208 Please refer to section 4.2 for primitive API specifications for individual problems.

209

210 3.2 EVOLUTIONARY SEARCH

In evolutionary program search, we typically start with a large population of random programs.
These programs are then sampled based on their fitness as parents. The parent programs create
new programs through crossover and mutation, resulting in a new population. Newer generations
improve over the previous as the population is getting optimized for the fitness function. The metric *M* can be used as the fitness function in our work.

Algorithm 1: DiSciPLE's learning loop **Input:** Observation set $\mathcal{D} = \{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$, metric \mathcal{M} , an objective prompt p_o , a set of primitive function $\mathcal{F} = \{f_1, f_2 \dots f_k\}$ **Hyperparameters:** Mutation probability ρ_m , total number of generations T, population size M, crossover prompt p_c , and mutation prompts p_m . **Output:** A hypothesis h^* in the form of a program that explains the observations best. $H^0 \leftarrow \{\} //$ Initialize a population of programs ² for i = 1, ..., M do $| h_i^0 \leftarrow \mathcal{LLM}(p_o)$ $H^0 \leftarrow H^0 \cup \{h_i^0\}$ $h^* \leftarrow h_i^0$ // Learn a \mathcal{D} that discriminates the training set $\{(x,y)\}$ 6 for t = 0, ..., T do $H^{t+1} \leftarrow \{\}$ for i = 1, ..., M do $h_{k_1}^t, h_{k_2}^t \leftarrow \mathsf{sample_parents}(H_t) \mathrel{//} \mathsf{Sample} \texttt{ parents} \texttt{ for crossover}$ $h_i^{t+1} \leftarrow \mathcal{LLM}(h_{k_1}^t, h_{k_2}^t, s(h_{k_1}^t; \mathcal{D}), s(h_{k_2}^t; \mathcal{D}), p_o, p_c) / / \text{ crossover operation}$ // critique and simplification $h_i^{t+1} \leftarrow \operatorname{critic}(h_i^{t+1}) h_i^{t+1} \leftarrow \operatorname{simplifier}(h_i^{t+1})$ $\stackrel{-}{H^{t+1}} \leftarrow H^{t+1} \ \cup \ \{h_i^{t+1}\}$ 17 return h^* ;

In our method, we keep the overall algorithm the same but replace key steps in using LLMs. First, at the start of the process, we provide the LLM with a prompt for the objective to generate the initial programs. To leverage the prior knowledge of LLMs, we use a prompt p_o that mentions the expert specified description of the quantity of interest: "Given a satellite image, write a function to estimate $\langle descr \rangle$ ". While an LLM cannot answer such a difficult scientific question without leveraging the observations \mathcal{D} , the prompt prevents the evolutionary algorithm from searching in completely random directions. As a result, our initial population is not entirely random.

Second, rather than using the symbolic methods of crossover and mutation, we use LLM to perform these operations. LLMs have common sense about programming and result in much better program modifications when performing crossover and mutations. More specifically, let $h_{k_1}^t$ and $h_{k_2}^t$ be two programs sampled from the t^{th} generation selected as parents based on the fitness function. To perform a crossover operation we pass, the objective prompt p_o , the two programs $h_{k_1}^t$ and $h_{k_2}^t$, their corresponding scores (using eq. (1)), along with a crossover prompt p_c to obtain a new program:

$$h_{k}^{t+1} = \mathcal{LLM}(h_{k_{1}}^{t}, h_{k_{2}}^{t}, s(h_{k_{1}}^{t}; \mathcal{D}), s(h_{k_{2}}^{t}; \mathcal{D}), p_{o}, p_{c})$$
(2)

The crossover prompt instructs the LLM to make use of the two-parent program and come up with a new program. The LLM is able to combine elements from the parents to produce something new as can be seen in fig. 2. Please refer to the appendix C.1 for more examples.

Similar to crossover, we also mutate a program with some probability using a mutation prompt p_m .

$${}^{m}h_{k}^{t+1} = \mathcal{LLM}(h_{k}^{t+1}, s(h_{k}^{t+1}; \mathcal{D}), p_{o}, p_{m})$$

$$\tag{3}$$

We present the exact input fed to the LLM for crossover and mutation in the appendix A. Note that both crossover and mutation operations also include the objective prompt preventing the LLM from generating programs far away from the objective.

2703.3FEATURE SET PREDICTION271

Hypotheses for many problems require combination of multiple feature. Optimizing for the correct combination of these features is challenging for an LLM, as we are not doing gradient-based optimization. Therefore instead of prompting the LLM to directly generate a predictor for *y*, we prompt is to create a list of predictive features. We can learn a linear regressor on top of this list and use the regressor with the list of features as hypotheses. Both the program and regressor are interpretable, therefore our hypothesis is also interpretable.

278 279 3.4 PROGRAM CRITIC

The only form of supervision our method gets is through the metric score s(h; D) created by evaluating the program h on the observations. However, since we perform crossover and mutation through a language model, we can provide finer-grained information to LLM in order to aid the search.

More specifically, we propose a critic that performs a finer-grained evaluation of the program that we get after crossover/mutation. Our critic performs a *stratified evaluation* by partitioning the observation data into multiple categories and evaluating the program on individual strata.

$$\mathcal{D} = d_1 \cup d_2 \cup \dots \cup d_c, \quad \text{where } d_i \cap d_j = \emptyset \text{ for } i \neq j \tag{4}$$

In all the applications where we use satellite images, we partition the observation dataset into landuse categories, using an open-world segmentation model. The critic obtains per-partition score $s(h; d_i)$ and prompts the LLM to improve the program on categories the model is bad. The addition of a critic improves the programs on data overlooked by programs, resulting in reliable programs.

3.5 PROGRAM SIMPLIFICATION

295 Successive steps of crossovers and mutations of programs result in large programs with many redun-296 dancies, hurting interpretability. We propose an analytical approach to simplify the programs and 297 remove the redundant parts of it. Our hypothesis are programs without loop, they can be represented 298 as a directed acyclic graphs (DAG) (we use the abstract syntax tree (AST); see fig. 5). In these 299 DAGs, all the constants and the arguments of the function are root nodes. Only the return statement 300 and the unused variables are the leaf nodes. Any leaf node that is not a return statement, is a piece 301 of code that is not needed and can be removed. We then recursively remove all the leaf nodes that are not return statements. 302

303 While removing such nodes (unreachable by the return statement) is useful, there could still be fea-304 tures that are returned by the program but are not contributing. Recall that we use linear regression 305 on the list of features returned by the program. The weights assigned to individual features by the 306 regression model are useful indicators of which features are redundant. We remove the features that 307 have a significantly smaller weight compared to the largest weight in the regression (a threshold of 308 5% works well). Removing these features from the return statement results in numerous newly created leaf nodes. We, therefore, redo the recursive leaf node removal to further simplify the program. 309 Each program generated by mutations and crossover is first improved through the critic and then 310 simplified before adding back to the population. Algorithm 1 shows the complete algorithm. 311

312 313

314

287 288 289

290

291

292 293

4 Results

315 4.1 IMPLEMENTATION DETAILS

We used the open-source LLM *llama-3-8b-instruct* (Dubey et al., 2024) served using the vLLM library (Kwon et al., 2023). A majority of our evaluation tasks use satellite images, so to allow inferring semantic information from it, we use a black-box open-world foundational model for satellite images, GRAFT (Mall et al., 2023). Some experiments use ground-truth annotations from Open-StreetMaps (Vargas-Munoz et al., 2020) as an alternative to disentangle the effect of segmentation.

We run our evolutionary method for T = 15 generations with a population size of M = 100. For all the problems, the input observation data comes from different geographical locations around the world. We split this data into three parts. Two-thirds of the easternmost observations are used



Figure 3: The best performing hypotheses for each of the 4 problems as Python programs (left in each card) and the corresponding DAG representation on the right. The DAG representation allows better visualization of the importance of different components. The thickness of the red edges determine how important that component is. A black edge represents computation; when removed it is either the same as one of its subsequent edges or removing it could result in a bug.

347

348 349

350 351

352

340

341

342

343

to create a training-testing split. The remaining one-third of the data is use to evaluate reliability (out-of-distribution generalization). We will release all four datasets for future research in this area.

4.2 SCIENTIFIC DOMAINS AND PROBLEMS

We apply DiSciPLE to two different problems in *Demography*: population density and poverty prediction. and to two problems in *Climate Science*: for AGB estimation and CSIF forecasting. In the following subsections, we present the observation datasets, metrics, and overview of primitives.

368

370

4.2.1 POPULATION DENSITY

Observation Dataset: The problem seeks to predict the population density by observing the satellite images of a region (Metzger et al., 2022; 2024). We obtain the population density values (y_i) for various locations in the USA by using ACS Community Surveys 5-year estimates (United States Census Bureau, 2024). Input observations (x_i) are sentinel-2 satellite images at a resolution of 10m (Drusch et al., 2012). For this experiment, we also use OpenStreetMaps masks (Vargas-Munoz et al., 2020) for 42 different land-use concepts (see appendix F.1) as part of the input.

Metric and Primitives: Population density values are aggregated at the county block group level.
 The predicted population densities are therefore also aggregated at the county block group level.
 The metric is the per-block group level average L2 error after applying a log transformation. Along
 with the arithmetic, and logical primitives (appendix B), we use open-vocabulary segmentation as a
 primitive. The segmentation function returns a binary mask for an input concept.

369 4.2.2 POVERTY INDICATOR

Observation Dataset: For poverty estimation, we use data from SustainBench (Yeh et al., 2021).
 The dataset contains coordinate location as input and wealth asset index as output.

Metric and Primitives: We use L2 error for each location as the evaluation metric. To obtain
 semantic land use information about a location, we first define a *get_satellite_image* function, that
 returns a sentinel-2 satellite image for any location. This can be used in conjunction with the open world satellite image recognition model to obtain semantic information about the world. Other than
 this we also include as primitives functions that return average annual temperature, precipitation,
 nightlight intensity, and elevation at the input location.

		Spatial					Temporal	
	Population Density		Poverty		AGB		CSIF	
	L2-Log	L1-Log	L1	RMSE	L1	RMSE	L1	RMSE
Mean	0.6696	0.6540	1.613	1.836	42.15	50.65	0.2521	0.3050
Concept Bottleneck	0.8298	0.7279	1.229	1.476	26.33	33.49	0.1474	0.1835
Deep Model - Small	0.4431	0.5006	1.238	1.637	30.72	37.03	0.0503	0.0727
Deep Model - Large	0.3974	0.4843	1.170	1.478	21.15	27.86	0.1487	0.1895
Zero-shot	0.4702	0.5371	1.525	1.754	38.80	46.41	0.2134	0.2610
Ours	0.2607	0.3778	1.077	1.314	24.79	32.99	0.0902	0.1260

378 Table 1: Performance of our hypotheses on unseen in-distribution observations across various prob-379 lems. In most cases, DiSciPLE produces better hypotheses (red is best and *blue* is second best). 380

4.2.3 ABOVEGROUND BIOMASS

Observation Dataset: Similar to poverty estimation, the observation variables are an input location and the output AGB estimate. We use NASA's GEDI (Dubayah et al., 2020) to obtain the observation value for three states in USA. We use data from the state of Massachusetts and Maine (North-East) as the train/test set and Washington (NorthWest) as the out-of-distibution set.

Metric and Primitives: We use L2 error as the metric and the same primitives as poverty estimation.

397 398 399

391 392

393

394

395

396

381 382

4.2.4 CONTIGUOUS SOLAR INDUCED CHLOROPHYLL FLUORESCENCE (CSIF)

400 Observation Dataset: The goal is to forecast CSIF values, by observing past CSIF and environmen-401 tal values. Unlike other problems, which could make use of spatial information and thus required 402 satellite images, CSIF forecasting is done without using satellite images. On the other hand, since 403 this is a forecasting problem, is uses past CSIF values as well as past and current values of environ-404 mental variable. The observation variables are an input location and output CSIF values. The data 405 comes from ERA5 climate data store (Hersbach et al., 2020).

406 Metric and Primitives: We use L2 error for each location as the evaluation metric. Along with the 407 mathematical and logical operations. The primitive allows obtaining present environmental variables 408 as well as a time series of past environmental variables such as minimum and maximum temperature 409 of past months, or soil moisture index (see appendix B). 410

4.3 EXPERIMENTAL SETUP 412

413 For the same set of training data we compare our best generated hypotheses with a set of baselines.

411

414

415 416

417

418

419

420

421

422

- 1. **Mean:** A naive baseline that use the mean of the training observation as the prediction.
- 2. Concept Bottleneck (Linear): Similar to Koh et al. (2020), we first extract a list of relevant features and train a linear classifier on it. This method is interpretable due to the bottleneck, however it is not very expressive (see appendix E.1).
 - 3. Deep models: We use deep models such as Resnets and LSTM as baseline (see appendix E.2 for details). We use a small and large variant for each.
 - 4. Zero-shot: This baseline tests how good would LLMs be on their own in generating hypotheses solely relying on prior knowledge without any observation. Since the generated programs can vary drastically, we report an average of 5 different zero-shot programs.
- 423 424 425 426
- 4.4 **RESULTS AND DISCUSSION**

427 We first test our hypothesis on unseen but *in-domain* observations/regions close to the regions used 428 for training (table 1). We observe that DiSciPLE outperforms all interpretable baselines. It can 429 even outperform a deep model in many cases, specifically on population density estimation, while being significantly more interpretable. DiSciPLE also outperforms zero-shot program inference 430 from LLMs suggesting that LLMs only have incomplete information about scientific domains and 431 our evolutionary process is necessary to identify a better hypothesis.

	Populatic	n Density	Spatial Poverty			GB	Temporal CSIF	
	L2-Log	L1-Log	L1	RMSE	L1	RMSE	L1	RMSE
Mean	0.6734	0.6561	1.591	1.844	74.15	83.02	0.2393	0.3018
Concept Bottleneck	0.7951	0.7112	1.257	1.504	44.19	63.52	0.1565	0.2029
Deep Model - Small	0.6623	0.5967	1.284	1.654	35.27	53.06	0.1061	0.1614
Deep Model - Large	0.4460	0.5115	1.344	1.741	35.41	70.30	0.1815	0.2435
Zero-shot	0.7020	0.6412	1.510	1.773	55.11	64.32	0.2190	0.2814
Ours	0.3807	0.4426	1.134	1.420	31.10	42.93	0.0882	0.1256

Table 2: Performance of our hypotheses on out-of-distribution observations across various problems.
 This shows the reliability of programs produced by DiSciPLE (red is best and *blue* is second best).

Are our hypotheses reliable? If a hypothesis is reliable it should be able to generalize to other regions. table 2 shows DiSciPLE to these baselines on such a out-of-distribution set. Here our approach outperforms all baselines *including deep networks*, suggesting that due to its interpretable-by-design representation our method learns model that can generalize better and overfit less.

Are our hypotheses data-efficient? Our methods are only trained on a maximum of 4000 observations. section 4.4 further shows that even when the amount of training data is reduced, our approach shows minimal degradation in performance compared to deep networks. This suggests that while deep models can learn to generalize with a lot more data, our model does not need as much data to begin with, making it data-efficient.



Figure 4: Performance of DiSciPLE compared to deep baselines as we reduce the amount of training observation (in terms of L2 error). The Oracle (blue) uses program learned from all observations, but uses only partial observation for parameter training. DiSciPLE (orange) uses partial observation during evolution as well. While the errors get worse as we reduce the observation data, the drop is significantly less severe for DiSciPLE compared to deep models, that tend to overfit.

Are our hypotheses interpretable? Our hypotheses are interpretable-by-design as we can visualize the factors contributing to performance. Fig. 3 shows such programs (left in each card) for all four of our problems. An expert who is working with our method to figure out such hypotheses can add/edit parts of the formula and figure out which/how much do each of these components matter.

We perform this step of understanding the influence of individual operation by removing each operation in our program, and measure its affects on the final score. The DAGs on the right of each program show the program structure and the red edges show the influence of each component proportional to the width. This visualization can allow experts to understand which which operations are important for the model. For example, in the program graph for population density fig. 3, we can see that for semantic concepts such as "highway" and "residential building" are very important.

Can our method perform better than expert humans? Our method would only be useful in real-world scenarios, if it can come up with stronger or comparable hypotheses to human experts. We test this on the task of AGB, by providing an expert (a PhD student actively working on AGB) with a user interface with the same information as our method. The experts took about 1.5 hours to use their domain knowledge and iterate over their program for AGB estimation. However, the best program they could come up with had an L1 error of 37.65 on the in-distribution set and 53.20 on the OOD set (compared to 24.79 and 31.10 for DiSciPLE). We figure this is primarily because

Table 3: Performance of our method as we successively remove the components. Both critic and simplification lead to performance improvement for our method.

			Test		00	DD
Feature-set pred.	Critic	Simplification	L2 log	L1 log	L2 log	L1 log
×	X	×	0.3159	0.4296	0.4835	0.5178
✓	X	×	0.2906	0.4049	0.4258	0.4826
1	1	×	0.2873	0.3984	0.4184	0.4684
1	1	1	0.2607	0.3778	0.3807	0.4426

experts need to spend more time on the problem. In general, experts would spend numerous days to come up with a good hypothesis, while our method can come up with a better hypothesis faster.

4.5 Ablations

How important is the role of feature-set prediction, critic, and simplification? Table 3 mea-501 sures the performance of our model on the task of population density as we successively add these 502 components to the evolutionary algorithm. The addition of feature set prediction instead of a single 503 feature helps, as it allows our method to learn expressive linear regression parameters instead of let-504 ting the LLM come up with them. Further adding critic results further improvement as the programs 505 start covering nicher concepts resulting in better unseen and OOD generalization. Finally adding 506 in simplification also improves hypothesis. We posit that simplification removes irrelevant features 507 preventing the LLM to focus on them when performing crossovers. 508

509 How important are common sense and prior knowledge of LLMs? The two major differences 510 or our method from symbolic regression are: 1) better crossover and mutation as LLMs can under-511 stand the meaning of the primitive functions. 2) use of prior knowledge for better guided search. 512 Therefore we remove these two sources of information and test how well can our method perform. To remove the understanding of functions we rename them with meaningless terms and remove the 513 descriptions. To remove the context of the problem we remove objective prompt. Without common 514 sense, the search cannot even progress away from the initial random programs, resulting in worse-515 than-mean results (L1 error of 0.84 vs 0.26 for DiSciPLE on density estimation). This suggests that 516 symbolic regression models, that have no understanding of open-world primitives, would struggle 517 to search. If we just remove the context of the problem the model does slightly better and can obtain 518 results better than the mean and zero-shot hypotheses (L1 error of 0.45). This suggests that while the 519 search is moving in the objective's direction, it is slow. (see appendix D.3 for the complete table). 520 We also present more ablations of our method in appendix D.

521 522 523

497

498 499

500

5 DISCUSSION AND CONCLUSION

524 **Limitations:** A limitation of DiSciPLE is that we can only differentiably optimize learnable pa-525 rameters in the last computational layer. This could miss out on hypotheses with useful parameters 526 in some intermediate computation layers. We attempted to make the whole pipeline differentiable, 527 however the model performance did not improve much. Many of the operations in our pipeline even 528 though differentiable have zero-gradient in large part of input space, making gradient optimization 529 challenging. Moreover, a completely differentiable hypothesis is even slower to optimize result-530 ing in much slower evolution. In future work, we plan to use initialization tricks for non-linear 531 optimization and second-order optimization to obtain even more expressive models.

532

Conclusion: We present DiSciPLE – an evolutionary algorithm that leverages the prior-knowledge and common sense abilities of LLMs to create *interpretable, reliable and data-efficient* hypotheses for real-world scientific problem. This allows us to create hypotheses that are more powerful than existing interpretable counterparts and more insightful than deeper uninterpretable models while being on-par in terms of performance. We shows its prowess on 2 different expert scientific domain on 4 different problems. We believe that using DiSciPLE in tandem with a human expert can rapidly speed up the scientific process and result in numerous novel discoveries.

540 REFERENCES 541

571

585

586

587

589

- Zeynep Akata, Scott Reed, Daniel Walter, Honglak Lee, and Bernt Schiele. Evaluation of output 542 embeddings for fine-grained image classification. In CVPR, 2015. 543
- 544 Aditya Chattopadhyay, Kwan Ho Ryan Chan, and Rene Vidal. Bootstrapping variational information pursuit with large language and vision models for interpretable image classification. In The 546 Twelfth International Conference on Learning Representations, 2024a. 547
- 548 Aditya Chattopadhyay, Ryan Pilgrim, and Rene Vidal. Information maximization perspective of orthogonal matching pursuit with applications to explainable ai. Advances in Neural Information 549 Processing Systems, 36, 2024b. 550
- 551 Mia Chiquier, Utkarsh Mall, and Carl Vondrick. Evolving interpretable visual classifiers with large 552 language models. 2024. 553
- 554 Miles Cranmer. Interpretable machine learning for science with pysr and symbolic regression. jl. 555 arXiv preprint arXiv:2305.01582, 2023.
- 556 Honghua Dong, Jiayuan Mao, Tian Lin, Chong Wang, Lihong Li, and Denny Zhou. Neural logic 557 machines. In International Conference on Learning Representations. 558
- 559 Matthias Drusch, Umberto Del Bello, Sébastien Carlier, Olivier Colin, Veronica Fernandez, Ferran 560 Gascon, Bianca Hoersch, Claudia Isola, Paolo Laberinti, Philippe Martimort, et al. Sentinel-2: 561 Esa's optical high-resolution mission for gmes operational services. *Remote sensing of Environ*-562 ment, 120:25-36, 2012.
- 563 Ralph Dubayah, James Bryan Blair, Scott Goetz, Lola Fatoyinbo, Matthew Hansen, Sean Healey, Michelle Hofton, George Hurtt, James Kellner, Scott Luthcke, et al. The global ecosystem dy-565 namics investigation: High-resolution laser ranging of the earth's forests and topography. Science 566 of remote sensing, 1:100002, 2020. 567
- 568 Abhimanyu Dubey, Abhinav Jauhri, Abhinav Pandey, Abhishek Kadian, Ahmad Al-Dahle, Aiesha 569 Letman, Akhil Mathur, Alan Schelten, Amy Yang, Angela Fan, et al. The llama 3 herd of models. 570 arXiv preprint arXiv:2407.21783, 2024.
- Kevin Ellis, Catherine Wong, Maxwell Nye, Mathias Sablé-Meyer, Lucas Morales, Luke Hewitt, 572 Luc Cary, Armando Solar-Lezama, and Joshua B Tenenbaum. Dreamcoder: Bootstrapping in-573 ductive program synthesis with wake-sleep library learning. In Proceedings of the 42nd acm 574 sigplan international conference on programming language design and implementation, pp. 835– 575 850, 2021. 576
- 577 Vittorio Ferrari and Andrew Zisserman. Learning visual attributes. Advances in neural information processing systems, 20, 2007. 578
- 579 Arya Grayeli, Atharva Sehgal, Omar Costilla-Reyes, Miles Cranmer, and Swarat Chaudhuri. Sym-580 bolic regression with a learned concept library. arXiv preprint arXiv:2409.09359, 2024. 581
- 582 Tanmay Gupta and Aniruddha Kembhavi. Visual programming: Compositional visual reasoning 583 without training. 2023 ieee. In CVF Conference on Computer Vision and Pattern Recognition 584 (CVPR), pp. 14953–14962, 2022.
 - Songhao Han, Le Zhuo, Yue Liao, and Si Liu. Llms as visual explainers: Advancing image classification with evolving visual descriptions. arXiv preprint arXiv:2311.11904, 2023.
- 588 Kaiming He, Xiangyu Zhang, Shaoqing Ren, and Jian Sun. Deep residual learning for image recognition. In Proceedings of the IEEE conference on computer vision and pattern recognition, pp. 590 770–778, 2016.
- Hans Hersbach, Bill Bell, Paul Berrisford, Shoji Hirahara, András Horányi, Joaquín Muñoz-Sabater, 592 Julien Nicolas, Carole Peubey, Raluca Radu, Dinand Schepers, et al. The era5 global reanalysis. Quarterly Journal of the Royal Meteorological Society, 146(730):1999–2049, 2020.

- Sepp Hochreiter and Jürgen Schmidhuber. Long short-term memory. *Neural computation*, 9(8): 1735–1780, 1997.
- Shaoli Huang, Zhe Xu, Dacheng Tao, and Ya Zhang. Part-stacked cnn for fine-grained visual cate-gorization. In *CVPR*, 2016.
- Justin Johnson, Bharath Hariharan, Laurens Van Der Maaten, Judy Hoffman, Li Fei-Fei,
 C Lawrence Zitnick, and Ross Girshick. Inferring and executing programs for visual reason ing. In *Proceedings of the IEEE international conference on computer vision*, pp. 2989–2998,
 2017.
- Elyor Kodirov, Tao Xiang, and Shaogang Gong. Semantic autoencoder for zero-shot learning. In
 CVPR, 2017.
- Pang Wei Koh, Thao Nguyen, Yew Siang Tang, Stephen Mussmann, Emma Pierson, Been Kim, and
 Percy Liang. Concept bottleneck models. 2020.
- Woosuk Kwon, Zhuohan Li, Siyuan Zhuang, Ying Sheng, Lianmin Zheng, Cody Hao Yu, Joseph Gonzalez, Hao Zhang, and Ion Stoica. Efficient memory management for large language model serving with pagedattention. In *Proceedings of the 29th Symposium on Operating Systems Principles*, pp. 611–626, 2023.
- Christoph H Lampert, Hannes Nickisch, and Stefan Harmeling. Attribute-based classification for zero-shot visual object categorization. *CoRR*, 2013.
- Jacky Liang, Wenlong Huang, Fei Xia, Peng Xu, Karol Hausman, Brian Ichter, Pete Florence, and
 Andy Zeng. Code as policies: Language model programs for embodied control. In 2023 IEEE
 International Conference on Robotics and Automation (ICRA), pp. 9493–9500. IEEE, 2023.
- ⁶¹⁸ Nour Makke and Sanjay Chawla. Interpretable scientific discovery with symbolic regression: a review. *Artificial Intelligence Review*, 57(1):2, 2024.
- 620
 621 Utkarsh Mall, Cheng Perng Phoo, Meilin Kelsey Liu, Carl Vondrick, Bharath Hariharan, and Kavita
 622 Bala. Remote sensing vision-language foundation models without annotations via ground remote
 623 alignment. arXiv preprint arXiv:2312.06960, 2023.
- Jiayuan Mao, Chuang Gan, Pushmeet Kohli, Joshua B Tenenbaum, and Jiajun Wu. The neuro symbolic concept learner: Interpreting scenes, words, and sentences from natural supervision.
 arXiv preprint arXiv:1904.12584, 2019.
- Sachit Menon and Carl Vondrick. Visual classification via description from large language models. *ICLR*, 2023.
- Nando Metzger, John E Vargas-Muñoz, Rodrigo C Daudt, Benjamin Kellenberger, Thao Ton-That
 Whelan, Ferda Ofli, Muhammad Imran, Konrad Schindler, and Devis Tuia. Fine-grained population mapping from coarse census counts and open geodata. *Scientific Reports*, 12(1):20085, 2022.
- Nando Metzger, Rodrigo Caye Daudt, Devis Tuia, and Konrad Schindler. High-resolution popula tion maps derived from sentinel-1 and sentinel-2. *Remote Sensing of Environment*, 314:114383, 2024.
- Juan Nathaniel, Gabrielle Nyirjesy, Campbell D Watson, Conrad M Albrecht, and Levente J Klein.
 Above ground carbon biomass estimate with physics-informed deep network. In *IGARSS 2023-2023 IEEE International Geoscience and Remote Sensing Symposium*, pp. 1297–1300. IEEE, 2023.
- Sarah Pratt, Ian Covert, Rosanne Liu, and Ali Farhadi. What does a platypus look like? generating customized prompts for zero-shot image classification. In *Proceedings of the IEEE/CVF International Conference on Computer Vision*, pp. 15691–15701, 2023.

Bernardino Romera-Paredes, Mohammadamin Barekatain, Alexander Novikov, Matej Balog,
M Pawan Kumar, Emilien Dupont, Francisco JR Ruiz, Jordan S Ellenberg, Pengming Wang,
Omar Fawzi, et al. Mathematical discoveries from program search with large language models. *Nature*, 625(7995):468–475, 2024.

648 649 650 651	Olaf Ronneberger, Philipp Fischer, and Thomas Brox. U-net: Convolutional networks for biomed- ical image segmentation. In <i>Medical image computing and computer-assisted intervention</i> - <i>MICCAI 2015: 18th international conference, Munich, Germany, October 5-9, 2015, proceed-</i> <i>ings, part III 18</i> , pp. 234–241. Springer, 2015.
652 653 654	Ozan Sener and Silvio Savarese. Active learning for convolutional neural networks: A core-set approach. In <i>International Conference on Learning Representations (ICLR)</i> , 2018.
655 656	Dídac Surís, Sachit Menon, and Carl Vondrick. Vipergpt: Visual inference via python execution for reasoning. <i>ICCV</i> , 2023.
657 658 659	Luming Tang, Davis Wertheimer, and Bharath Hariharan. Revisiting pose-normalization for fine- grained few-shot recognition. In <i>CVPR</i> , 2020.
660 661	Qwen Team. Qwen2.5: A party of foundation models, September 2024. URL https://qwenlm.github.io/blog/qwen2.5/.
662 663 664	United States Census Bureau. American community survey 5-year estimates. https://www.census.gov/programs-surveys/acs, 2024. Accessed: 2024-10-01.
665 666 667 668	Grant Van Horn, Steve Branson, Ryan Farrell, Scott Haber, Jessie Barry, Panos Ipeirotis, Pietro Perona, and Serge Belongie. Building a bird recognition app and large scale dataset with citizen scientists: The fine print in fine-grained dataset collection. In <i>Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition</i> , pp. 595–604, Boston, MA, 2015.
669 670 671 672	John E Vargas-Munoz, Shivangi Srivastava, Devis Tuia, and Alexandre X Falcao. Openstreetmap: Challenges and opportunities in machine learning and remote sensing. <i>IEEE Geoscience and</i> <i>Remote Sensing Magazine</i> , 9(1):184–199, 2020.
673 674	Michael Xie. <i>MAPPING POVERTY WITH SATELLITE IMAGERY</i> . PhD thesis, STANFORD UNI-VERSITY, 2017.
675 676 677 678	Christopher Yeh, Chenlin Meng, Sherrie Wang, Anne Driscoll, Erik Rozi, Patrick Liu, Jihyeon Lee, Marshall Burke, David B Lobell, and Stefano Ermon. Sustainbench: Benchmarks for monitor- ing the sustainable development goals with machine learning. <i>arXiv preprint arXiv:2111.04724</i> , 2021.
679 680 681	Xixian Yong and Xiao Zhou. MuseCL: Predicting urban socioeconomic indicators via multi- semantic contrastive learning. <i>CoRR</i> , 2024.
682 683 684	Yao Zhang, Joanna Joiner, Seyed Hamed Alemohammad, Sha Zhou, and Pierre Gentine. A global spatially contiguous solar-induced fluorescence (csif) dataset using neural networks. <i>Biogeosciences</i> , 15(19):5779–5800, 2018.
685 686 687	Bolei Zhou, Yiyou Sun, David Bau, and Antonio Torralba. Interpretable basis decomposition for visual explanation. In <i>ECCV</i> , 2018.
689	
690	
691	
692	
693	
694	
695	
696	
697	
698	
699	
700	
/ 11 1	

702 **CROSSOVER AND MUTATION PROMPTS** А 703

704

```
For crossover we take two programs and their corresponding scores. We use the following prompt
705
       for crossover:
706
       \{\{h_1\}\}\
707
708
       This program has a score of \{\{s(h_1)\}\}.
709
       \{\{h_2\}\}\
710
       This program has a score of \{\{s(h_2)\}\}.
711
712
       Can you write a function that gives higher score? Feel free to
713
       combine elements that worked from both programs. Only give me
714
       code.
715
       And similarly we use the following prompt for random mutation.
716
717
       \{\{h\}\}\
718
       This program has a score of \{\{s(h)\}\}.
719
       Can you edit this code to write a better function for the problem?
720
       Only give me code.
721
722
723
           PROBLEM SPECIFIC PRIMITIVE DESCRIPTION
       B
724
725
       B.1 PRIMITIVES AND THEIR DESCRIPTIONS FOR POPULATION DENSITY
726
727
728 1
729 <sup>2</sup>
        def elementwise_max(matrix1, matrix2):
             .....
    3
730 4
             Compute the element-wise maximum of two matrices.
731 5
732 6
             Parameters:
               matrix1 (numpy.ndarray): First input matrix.
733 7
                matrix2 (numpy.ndarray): Second input matrix.
734 8
735 <sub>10</sub>
             Returns:
736 11
                numpy.ndarray: Element-wise maximum of the input matrices.
             .....
737 12
738 13
739 <sup>14</sup>
        def elementwise_min(matrix1, matrix2):
739
740 15
16
             .....
             Compute the element-wise minimum of two matrices.
741 17
742 18
             Parameters:
                matrix1 (numpy.ndarray): First input matrix.
743 19
744<sup>20</sup>
                 matrix2 (numpy.ndarray): Second input matrix.
744
745 21
22
             Returns:
746 <sub>23</sub>
                 numpy.ndarray: Element-wise minimum of the input matrices.
             .....
747 24
748<sup>25</sup>
749 <sup>26</sup> <sub>27</sub>
        def elementwise_sum(matrix1, matrix2):
750 <sub>28</sub>
             Compute the element-wise sum of two matrices.
751 29
752 30
             Parameters:
                 matrix1 (numpy.ndarray): First input matrix.
753 <sup>31</sup>
754 <sup>32</sup>
                 matrix2 (numpy.ndarray): Second input matrix.
   33
755 34
             Returns:
```

```
756
                  numpy.ndarray: Element-wise sum of the input matrices.
757 35
              .....
758 36
37
759 38
         def elementwise_product(matrix1, matrix2):
760 39
              .....
              Compute the element-wise product of two matrices.
761 40
762 41
763 <sup>42</sup> <sub>43</sub>
              Parameters:
                  matrix1 (numpy.ndarray): First input matrix.
764 <sub>44</sub>
                  matrix2 (numpy.ndarray): Second input matrix.
765 45
766 46
             Returns:
767 <sup>47</sup>
                  numpy.ndarray: Element-wise product of the input matrices.
768 <sup>48</sup>
    49
769 50
        def elementwise_division(matrix1, matrix2):
770 51
              .....
             Compute the element-wise division of two matrices.
771 52
772 <sup>53</sup>
773 <sup>54</sup>
55
             Parameters:
                 matrix1 (numpy.ndarray): First input matrix.
774 56
                  matrix2 (numpy.ndarray): Second input matrix.
775 57
776 58
             Returns:
777 59
                 numpy.ndarray: Element-wise division of the input matrices.
              .....
778 <sup>60</sup><sub>61</sub>
779<sub>62</sub>
        def matrix_scalar_multiplication(matrix, scalar):
780 63
              Perform matrix scalar multiplication.
781 64
782<sup>65</sup>
             Parameters:
783 <sup>66</sup>
   67
                matrix (numpy.ndarray): Input matrix.
784 68
                  scalar (int or float): Scalar value.
785 69
             Returns:
786 70
              numpy.ndarray: Result of matrix scalar multiplication.
"""
787 71
788 72
73
789 <sub>74</sub>
         def elementwise_log(matrix):
790 75
              .....
             Compute the element-wise logarithm of a matrix.
791 76
792 <sup>77</sup>
793 78
79
              Parameters:
                 matrix (numpy.ndarray): Input matrix.
794 <sub>80</sub>
795 81
             Returns:
796 82
                  numpy.ndarray: Element-wise logarithm of the input matrix.
              .....
797<sup>83</sup>
798 <sup>84</sup> <sub>85</sub>
         def elementwise_exponentiate(matrix, base):
799 <sub>86</sub>
800 87
              Compute the exponentiation of each element of a matrix with a given
              \rightarrow base.
801
802 88
803 <sup>89</sup>
90
             Parameters:
                 matrix (numpy.ndarray): Input matrix.
804 <sub>91</sub>
                  base (int or float): Base value.
805 92
806 93
             Returns:
              numpy.ndarray: Exponentiated matrix.
807 <sup>94</sup>
808 95
   96
809 <sub>97</sub>
        def min_pixel_distance_to_mask(mask):
```

```
810
              .....
811 98
812<sup>99</sup>
100
              Compute the minimum pixel distance from each pixel to a mask.
813<sub>101</sub>
              Parameters:
814102
                 mask (numpy.ndarray): Binary mask array where 1 represents the
                  \rightarrow mask and 0 represents the background.
815
816<sup>103</sup>
817<sup>104</sup><sub>105</sub>
              Returns:
                 numpy.ndarray: Minimum pixel distance to the mask.
818<sub>106</sub>
              .....
819107
        def segment(im, text_prompt="trees"):
820 108
821 <sup>109</sup>
822<sup>110</sup>
              Segments a satellite image based on a text prompt. The text prompt
              → can only take one concept at a time.
823111
824112
              Parameters:
                  im (numpy.ndarray): An rgb satellite image.
825113
                  text_prompt (str): a text prompt
826114
827<sup>115</sup>
   116
              Returns:
828<sub>117</sub>
                 mask (numpy.ndarray): Binary mask array where 1 represents the
829
                   \rightarrow mask and 0 represents the background.
830118
             Example:
831 <sup>119</sup>
              text_prompt can be but not limited to these:
832<sup>120</sup>
              ["tennis", "skate park", "football field", "swimming pool",
              ↔ "cemetery", "multi-storey garage", "golf", "roundabout",
833
              ↔ "parking lot", "supermarket", "school", "marina", "baseball
              → field", "fall", "pond", "airport", "beach", "bridge", "religious
→ building", "residential building", "warehouse", "office
834
835
              ↔ building", "farmland", "university building", "forest", "lake",
836

    → "nature reserve", "park", "sand", "soccer field", "equestrian
    → club", "shooting range", "ice-rink", "commercial area",

837
838
                   "garden", "dam", "railroad", "highway", "river", "wetland",
              \hookrightarrow
839
              ↔ "non-residential buildings", "coastline"]
              .....
840121
841
842
843
        B.2 PRIMITIVES AND THEIR DESCRIPTIONS FOR AGB AND POVERTY PREDICTION
844
        The evolutionary search uses all the above defined functions, plus the following:
845
846
847 1
848 <sup>2</sup><sub>3</sub>
        def get_satellite_image(location):
849 4
              .....
             Get the satellite image for a given location.
850 5
             Parameters:
851 6
852 <sup>7</sup>
                  location (tuple): Tuple containing the latitude and longitude of
                  → the location.
853 <sub>8</sub>
             Returns:
854 9
                  numpy.ndarray: Satellite image for the location.
855 10
             Can be used for segmentation ONLY.
856 11
              .....
857 <sup>12</sup>
   13
858 14
         def get_temperature(location):
859 15
              .....
860 16
              Get the average annual temperature for a given location.
861 17
              Parameters:
862 18
                   location (tuple): Tuple containing the latitude and longitude of
                   \rightarrow the location.
863 19
              Returns:
```

```
864
                  float: Temperature for the location normalized between 0 and
865 20
866 21
                   ↔ 255.
              .....
867 <sub>22</sub>
868 23
        def get_precipitation(location):
              .....
869 24
             Get the average annual precipitation for a given location.
870 <sup>25</sup>
871 <sup>26</sup>
27
             Parameters:
                  location (tuple): Tuple containing the latitude and longitude of
872
                  → the location.
873 <sub>28</sub>
             Returns:
                  float: Precipitation for the location between 0 and 255.
874 29
              .....
875 <sup>30</sup>
876 <sup>31</sup>
    32
        def get_elevation(location):
877 33
             .....
878 34
             Get the elevation for a given location.
879 35
             Parameters:
                  location (tuple): Tuple containing the latitude and longitude of
880 36
881 <sub>37</sub>
                  \rightarrow the location.
             Returns:
882 <sub>38</sub>
                  float: Digital Elevation for the location (scaled 0-8000) to
883
                  ↔ 0-255.
             .....
884 39
885 <sup>40</sup>
886 41
        def get_nightlight_intensity(location):
              .....
887 43
             Get the average annual nightlight intensity for a given location.
888 44
             Parameters:
                  location (tuple): Tuple containing the latitude and longitude of
889 45
                  → the location.
890
891 <sup>46</sup>
             Returns:
   47
                  float: Nightlight intensity for the location (between 0 and 1).
892 <sub>48</sub>
893 49
        def get_average(segmented_image):
894 50
             .....
895 51
896 <sup>52</sup>
53
             Get the average pixel value of a segmented image.
897 54
             Parameters:
898 55
                 segmented_image (numpy.ndarray): Segmented image.
899 56
900 57
             Returns:
901 <sup>58</sup>
                  float: Average pixel value of the segmented image.
    59
902
903
904
        B.3 PRIMITIVES AND THEIR DESCRIPTIONS FOR CSIF FORECASTING
905
906
        For CSIF forecast, the API borrows mathematical and logical functions from above. Additionally it
907
        has the following to obtain more environmental variables.
908
909 1
910 2
        def get_historical_csif(locations, num_months=36):
911 <sup>3</sup>
              .....
912 4
             Get historical CSIF (contiguous solar induced chlorophyll
    5
913
              ↔ fluorescence) time-series data for the last given number of
914
              \hookrightarrow months.
915 6
             Parameters:
916 <sup>7</sup>
                  locations: a list of locations in a specific format.
    8
917
```

```
17
```

 \hookrightarrow

for.

num_months (int): Last number of months to get time-series data

```
918
              Returns:
919<sup>10</sup>
920 11
                 numpy.ndarray: Historical time series data for the given
                   → locations. size = (len(locations), num_months)
921 12
              .....
922 13
923 14
         def get_historical_min_temperature(locations, num_months=36):
              .....
924 15
925 <sup>16</sup>
              Get historical minimum temperature time-series data for the last
              \rightarrow given number of months.
926 <sub>17</sub>
927 18
              Parameters:
                   locations: a list of locations in a specific format.
928 19
                  num_months (int): Last number of months to get time-series data
929 <sup>20</sup>
930 <sub>21</sub>
                   \hookrightarrow for.
              Returns:
931 <sub>22</sub>
                 numpy.ndarray: Historical time series data for the given
932
                   → locations. size = (len(locations), num_months)
              .....
933 23
934 <sup>24</sup>
935 <sup>25</sup>
         def get_historical_max_temperature(locations, num_months=36):
    26
936 <sub>27</sub>
              .....
937 28
              Get historical maximum temperature time-series data for the last
              \rightarrow given number of months.
938
939 <sup>29</sup>
940 <sup>30</sup><sub>31</sub>
              Parameters:
                  locations: a list of locations in a specific format.
941 32
                   num_months (int): Last number of months to get time-series data
942
                  \rightarrow for.
              Returns:
943 33
                 numpy.ndarray: Historical time series data for the given
944 <sup>34</sup>
945 <sub>35</sub>
                  → locations. size = (len(locations), num_months)
              .....
946 <sub>36</sub>
947 37
         def get_historical_radiation(locations, num_months=36):
948 38
              .....
949<sup>39</sup>
950 <sup>40</sup>
              Get historical solar radiation time-series data for the last given
              \rightarrow number of months.
951 <sub>41</sub>
952 42
              Parameters:
                   locations: a list of locations in a specific format.
953 43
954 <sup>44</sup>
                   num_months (int): Last number of months to get time-series data
                   \rightarrow for.
955 45
              Returns:
956 <sub>46</sub>
                  numpy.ndarray: Historical time series data for the given
957
                   → locations. size = (len(locations), num_months)
              .....
958 47
959 <sup>48</sup>
960 <sup>49</sup>
50
         def get_historical_precipitation(locations, num_months=36):
              .....
961 <sub>51</sub>
              Get historical precipitation time-series data for the last given
962
              \hookrightarrow number of months.
963 52
              Parameters:
964 <sup>53</sup>
965 <sup>54</sup>
                  locations: a list of locations in a specific format.
                  num_months (int): Last number of months to get time-series data
    55
966
                   \rightarrow for.
967 56
              Returns:
                  numpy.ndarray: Historical time series data for the given
968 57
                   → locations. size = (len(locations), num_months)
969
              .....
970 <sup>58</sup>
    59
971 <sub>60</sub>
```

```
972
         def get_historical_photoperiod(locations, num_months=36):
973 61
              ......
974 <sup>62</sup>
975<sup>63</sup>
              Get historical photoperiod time-series data for the last given
              \rightarrow number of months.
976 64
977 65
              Parameters:
                   locations: a list of locations in a specific format.
978 <sup>66</sup>
979 <sup>67</sup>
                   num_months (int): Last number of months to get time-series data
                   \hookrightarrow for.
980 <sub>68</sub>
              Returns:
981 69
                  numpy.ndarray: Historical time series data for the given
                   → locations. size = (len(locations), num_months)
982
              .....
983 <sup>70</sup>
984 <sup>71</sup>
72
985 73
         def get_historical_swvl1(locations, num_months=36):
986 74
              Get historical soil water content in the first layer time-series
987 75
              \leftrightarrow data for the last given number of months.
988
989 <sup>76</sup>
77
              Parameters:
990 <sub>78</sub>
                   locations: a list of locations in a specific format.
991 79
                  num_months (int): Last number of months to get time-series data
                   \rightarrow for.
992
993 <sup>80</sup>
              Returns:
994 <sup>81</sup>
                  numpy.ndarray: Historical time series data for the given
                   → locations. size = (len(locations), num_months)
995 <sub>82</sub>
              .....
996 83
997 84
         def get_present_min_temperature(locations):
998<sup>85</sup>
999 <sup>86</sup>
    87
              Get present minimum temperature data for the given locations.
100088
100189
              Parameters:
                  locations: a list of locations in a specific format.
100290
              Returns:
1003<sup>91</sup>
                 numpy.ndarray: Present minimum temperature data for the given
1004<sup>92</sup>
                   → locations. size = (len(locations),)
1005<sub>93</sub>
              .....
100694
100795
         def get_present_max_temperature(locations):
1008<sup>96</sup>
1009<sup>97</sup>
98
              Get present maximum temperature data for the given locations.
1010<sub>99</sub>
              Parameters:
101100
                  locations: a list of locations in a specific format.
              Returns:
101<u>1</u>01
1013^{02}
                 numpy.ndarray: Present maximum temperature data for the given
                  → locations. size = (len(locations),)
1014
103
101504
101 {\textstyle \textcircled{}_{05}}
         def get_present_radiation(locations):
              .....
101706
              Get present solar radiation data for the given locations.
101807
1019^{108}_{109}
              Parameters:
1020<sub>10</sub>
                  locations: a list of locations in a specific format.
102111
              Returns:
                  numpy.ndarray: Present solar radiation data for the given
102912
                   → locations. size = (len(locations),)
1023
              .....
1024<sup>113</sup>
   114
1025<sub>15</sub>
         def get_present_precipitation(locations):
```

```
1026
             .....
102716
1028<sup>117</sup>
118
             Get present precipitation data for the given locations.
1029<sub>19</sub>
             Parameters:
103Q<sub>20</sub>
                 locations: a list of locations in a specific format.
103121
             Returns:
                 numpy.ndarray: Present precipitation data for the given
103<sup>422</sup>

→ locations. size = (len(locations),)

1033
123
             .....
1034<sub>124</sub>
103525
        def get_present_photoperiod(locations):
             .....
103£26
103727
             Get present photoperiod data for the given locations.
1038^{128}_{129}
             Parameters:
1039<sub>130</sub>
                 locations: a list of locations in a specific format.
104Q<sub>31</sub>
             Returns:
                 numpy.ndarray: Present photoperiod data for the given locations.
104132

→ size = (len(locations),)

1042
             .....
1042
1043
134
1044<sub>135</sub>
        def get_present_swvl1(locations):
104536
             ....
             Get present soil water content in the first layer data for the given
104637
             \rightarrow locations.
1047
1048
139
             Parameters:
1049<sub>140</sub>
              locations: a list of locations in a specific format.
105041
             Returns:
                numpy.ndarray: Present soil water content in the first layer
105142
                  → data for the given locations. size = (len(locations),)
1052
             .....
1053^{143}_{144}
1054
145
1055
1056
1057
           MORE EXAMPLES
1058
       С
1059
1060
       C.1 MORE EXAMPLES OF CROSSOVER
1061
       Parent 1
1062
1063
        def estimator(location):
1064 \frac{1}{2}
             images = get_satellite_image(location)
1065<sub>3</sub>
             temperature = get_temperature(location)
1066 4
             precipitation = get_precipitation(location)
             nightlight = get_nightlight_intensity(location)
1067 5
             return temperature, precipitation, elevation, nightlight
1068<sup>6</sup>
1069
1070
1071
       Parent 2
1072
        def estimator(location):
1073 1
             images = get_satellite_image(location)
1074<sup>2</sup>
1075<sup>3</sup><sub>4</sub>
             roads = segment(images, 'roads')
             avg_roads = get_average(roads)
1076 <sub>5</sub>
             avg_nightlight = get_nightlight_intensity(location)
1077
1078
1079
       Program after crossover
```

```
1080
        def estimator (location):
1081<sup>1</sup>
            images = get_satellite_image(location)
1082<sup>2</sup>
            roads = segment(images, 'roads')
1083 4
            avg_roads = get_average(roads)
1084 5
            avg_nightlight = get_nightlight_intensity(location)
            avg_precipitation = get_precipitation(location)
1085 6
            avg_temperature = get_temperature(location)
1086 7
            return (avg_roads * avg_precipitation * avg_nightlight *
1087 <sup>8</sup>
             \rightarrow avg_temperature)
1088
1089
1090
       Parent 1
1091
1092
        def estimator(location):
1093 2
            images = get_satellite_image(location)
1094<sub>3</sub>
            forests = segment(images, 'forests')
1095 4
            avg_forests = get_average(forests)
            roads = segment(images, 'roads')
1096 5
1097 <sup>6</sup>
            avg_roads = get_average(roads)
            elevation = get_elevation(location)
1098
            nightlight_intensity = get_nightlight_intensity(location)
1099 9
            return avg_forests, avg_roads, 1 / (1 + elevation), 1 / (1 +
1100
            → nightlight_intensity)
1101
1102
1103
       Parent 2
1104
1105
        def estimator(location):
1106 2
            images = get_satellite_image(location)
            poverty_mask = segment(images, 'poverty')
1107<sup>3</sup>
            return get_average(poverty_mask)
1108 4
1109
1110
       Program after crossover
1111
1112
1113 2
        def estimator(location):
            images = get_satellite_image(location)
1114 3
            poverty_mask = segment(images, 'poverty')
1115 4
1116 5
            education_mask = segment(images, 'education')
            health_mask = segment(images, 'health')
1117 7
            avg_poverty = get_average(poverty_mask)
1118 8
            avg_education = get_average(education_mask)
            avg_health = get_average(health_mask)
11199
1120<sup>10</sup>
            elevation = get_elevation(location)
            nightlight_intensity = get_nightlight_intensity(location)
   11
1121<sub>12</sub>
            return avg_poverty, avg_education, avg_health, 1 / (1 + elevation),
1122
             \rightarrow 1 / (1 + nightlight_intensity)
1123
1124
1125
       C.2 ILLUSTRATION OF SIMPLIFCATION
1126
1127
       For the following program, we show the steps of simplification. Program after crossover
1128
1129 1
        def estimator(im):
            building_mask = segment(im, "residential building")
1130 2
            nr_mask = segment(im, "non-residential buildings")
1131<sup>3</sup>
1132 <sup>4</sup>
             vegetation_mask = segment(im, "forest")
            water_mask = segment(im, "lake")
road_mask = segment(im, "highway")
1133_{6}^{-1}
```

Table 4:	different	LLMs
----------	-----------	------

			_		_		
		1116	T	est	00	D	
			L2 log	L1 log	L2 log	LI log	
		Qwen2.5-/b-instruct	0.2950	0.4039	0.4263	0.4716	
		llama-3 1-8h-instruct	0.2020	0.3842	0.3940	0.4518	
		llama-3.1-70b-instruct	0.2896	0.3958	0.4223	0.4663	
		Deep Model - Large	0.3974	0.4843	0.4460	0.5115	
		Deep moder Darge	010771	011012	011100	010110	
	building nr_dista vegetati water_di road_dis	_distance = min_pix nce = min_pixel_dis on_distance = min_p stance = min_pixel_ tance = min_pixel_c	kel_dist stance_t pixel_di _distance distance	ance_to o_mask(stance_ e_to_mas _to_mas	o_mask(k [nr_mas] to_mas] ask(wate sk(road_	ouilding_m c) c(vegetati er_mask) _mask)	ask) on_mask)
	return b ↔ road	uilding_distance, r _distance	nr_dista	nce, ve	egetatio	on_distanc	e,
Usin bui that	ng regression lding_dis and recursive	n weights our meth tance is not a useful ely all the leaf nodes.	nod figur value to	res out be return	that the d. So i	ne left mo in the third s	ost branc step we re
highwa	residential	non-residential im forest	lake	sidential	non-residential		forest
segn	ient segment	segmént segment	Segment	segment	segment	segment	Setlin
assign roa	ad_mask assign building_mas	k assign nr_mask assign vegetation_mask ass	ign water_mask 1	ussign building_mask	assign nr_mas	k assign vegetation_r	mask assign ro
distance_t	distance_transform	distance_transform distance_transform dist	tance_transform	distance_transform	distance_transfo	distance_transfor	rm distance
assign	road_distance assign building dista	ince assign nr_distance assign vegetation_distance assig	n water_distance	assign building_distance	assign nr_dista	ance assign vegetation_dista	ance assign road_o
			-			tipið	
		mbac.				Ţ	
		return					
000-000	dantial	En Event	birthumu	residential	non-resident		forest
non-resi							
seg	ment	segment	Segment	segment	segment	segment	
assign	nr_mask	assign vegetation_mask as	sign road_mask 🔰 a	ssign building_mask	assign nr_m	assign vegetation	n_mask assign
_	transform	distance_transform	tance_transform	distance_transform	distance_tran	sform distance_trans	sform

Figure 5: Process of simplification illustrated over a function.

1182 D ADDITIONAL ABLATIONS

1184 D.1 USING DIFFERENT LLMS

A key contribution of our work is to leverage the common-sense knowledge in LLMs to improve evolutionary search. So, it is natural to question whether (a) LLMs (with similar capacity) trained with a different large corpus of text would generate hypotheses with different levels of reliability and

1188 Table 5: Perfomance of our method when removing the context of the problem (objective prompt 1189 from the evolution, and when renaming and not describing the primitive functions to the LLM. We 1190 see significant drops in performance in both cases, suggesting that both the common-sense and prior knowledge of LLM is important to perform efficient evolutionary search.) 1191

1192			Mathad	Lilea	I 2 log		
1193			Method	LI log	L2 log		
1194			No common-sense	0.8401	0.7186		
1195			DiSciPLE full	0.4498	0.3778		
1196		•		0.2001	0.0110		
1197							
1198	(b) L1	Ms with larger capacity	would produce more	e reliable	hypothese	es. We answei	• these questions
1199	by tes	sting recent LLMs: Owe	n-2.5/7b (Team, 202	4), llama-	-3/8b, llan	na-3.1/8b, lla	ma-3.1/70b. For
1200	ease of	of experimentation, we re	educe the number of	generatio	ons to 10 a	and the popul	ation size to 60.
1201	We re	port the results in table	4. DiSciPLE works	robustly	with vario	ous LLMs and	d could generate
1202	more	reliable models/hypothe	ses than the Deep M	lodel. Wl	hile the pe	erformance of	f the hypotheses
1203	varies	s, we do not observe any	discernible differenc	e among	the variou	s hypotheses.	
1204							
1205	D.2	ALBATION ON NOISY/	UNRELIABLE PRIMI	TIVES			
1206	- ·						
1207	To in	vestigate how accurate/r	obust should the und	derlying t	black-box	model be?, y	we corroded the
1208	OSM	maps with a 3x3 convo	olution and ran Disc	$c_{1}PLE$ to	generate	a new hypot	nesis. With the
1209	perfo	rmance compared to clear	n OSM maps (I 2 log)	010.3/13	2626	t set — a larg	e degradation in
1210	peno	initialitie compared to ciea	ii OSWi iiiaps (L2 i0g	g c1101. 0.	.2020).		
1211	D 2	A DI ATIONS OF LLM	COMMON CENCE AN		KNOWL EI	DCE	
1212	D.5	ADLATIONS OF LLIVI (OMMON-SENSE AN	DPRIOR	KNUWLE	DGE	
1213	table	5 shows more extensive e	evaluation performing	g evaluati	on withou	t primitive un	derstanding and
1214	probl	em understanding.	1	C		1	e
1215	•	-					
1216	E I	EXPERIMENTAL SET	TUP				
1217							
1218	E.1	MORE DETAILS ON CO	NCEPT BOTTLENECI	K BASELI	NES		
1219	2.11				1120		
1220	For the	he CSIF task, the conce	pt bottleneck feature	es are ave	erage of p	oast CSIF and	d environmental
1221	variał	ole and the current enviro	onmental variable. F	or all the	other task	s, the bottlen	eck features are
1222	42 ca	tegories of segments obta	ained from either OS	M or GR	AFT and t	he environme	ental variable.
1223							
1224	E.2	MORE DETAILS ON DE	EP MODELS BASELI	NE			
1225	_						
1226	For sp	patial tasks, we use a Res	Net-18 (He et al., 201	6) based	U-Net (Ro	onneberger et	al., 2015) (Deep
1227	Mode	el-Large). Since our mod	fully convolution -1	data-effic	ient, to pi	event overfitt	ing, we also try
1228	a sina	Iner Dackbone of 4-layer	huber (1007) of diff	arent den	th as our	big (12 laver	c_{s} and small (4)
1229	mode	is mountaire a scilling		erent uep	ui as oul	oig (12-layel	s) and sman (4)
1230	mout	10					
1231	_						
1232	ΓI	MORE DETAILS					
1233							

1234 F.1 LIST OF 42 LAND-USE CONCEPTS

1235 Table 6 show the list of 42 concepts extracted from OpenStreetMaps and also used in GRAFT to get 1236 partitions for critic. 1237

- 1238
- 1239
- 1240
- 1241

1242						
1243						
1244						
1245						
1246						
1247						
1248						
1249						
1250						
1251						
1252						
1253						
1254						
1255						
1256						
1257						
1258						
1259						
1260						
1261						
1262						
1263						
1264						
1265	tannia aqueta	alrata manlr	amariaan faathall fald	autimmina naal	aamatami	nond
1266	golf course	roundabout	narking lot	swimming poor	school	marina
1267	baseball	waterfall	multi-storey parking garage	airport	beach	bridge
1268	religious building	residential building	university building	office	farmland	warehouse
1269	forest	lake	nature reserve	park	sandy area	soccer field
1270	equestrian center railroad	shooting range highway	non residential buildings river or stream	commercial area wetland	garden ice-rink	dam coastline
1271	Table 6: List of cor	conts extracted from	OSM and also used via CP	AFT for critic date	stratificatio	n
1273	Table 0. List of col	icepts extracted from	OSIM and also used via OK	AFT for critic data	stratificatio	11.
1274						
1275						
1276						
1277						
1278						
1279						
1280						
1281						
1282						
1283						
1284						
1285						
1286						
1287						
1288						
1289						
1290						
1291						
1292						
1000						
1293						
1293						
1293 1294 1295						