
Importance of Directional Feedback for LLM-based Optimizers

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 We study the potential of using large language models (LLMs) as an interactive
2 optimizer for solving maximization problems on a text space using natural language
3 and numerical feedback. Inspired by the classical optimization literature, we
4 classify the natural language feedback into directional and non-directional, where
5 the former is a generalization of the first-order feedback to the natural language
6 space. We find that LLMs are especially capable of optimization when they
7 are provided with directional feedback. Based on this insight, we design a new
8 LLM-based optimizer that synthesizes directional feedback from the historical
9 optimization trace to achieve reliable improvement over iterations. Empirically, we
10 show our LLM-based optimizer is more stable and efficient in solving optimization
11 problems, from maximizing mathematical functions to optimizing prompts for
12 writing poems, compared with existing techniques.

13 1 Introduction

14 Owing to their capability to produce a diverse range of outputs similar to those of humans, large
15 language models (LLMs) are a powerful component for solving many difficult problems involving
16 natural language, including planning [3], interacting with users, understanding documents [4], and
17 producing executable code [2]. In addition to harnessing LLMs in these *generative* roles, several
18 recent works have used LLMs for *optimization*. So far, these efforts, such as APO [7] and OPRO [11],
19 have focused on optimization of a very specific kind – employing LLMs to produce prompts that
20 improve (another) LLM’s performance. In this work, we argue that LLMs’ potential extends much
21 further, to general optimization problems. We showcase that LLMs are capable of optimizing
22 entities as dissimilar as mathematical functions and poems if they are provided with *directional*
23 *feedback*.

24 The notions of directional and non-directional feedback arise naturally in many interactive decision-
25 making domains and are tied to the classical optimization literature [1]. Typically, a numerical
26 optimizer iterates over two steps. The first step aims to identify a “search direction” for improvement.
27 This information is provided to the optimizer by an oracle, oftentimes a first-order oracle, and
28 can be viewed as directional feedback. The second step decides what to change about the input.
29 The applicability of various optimization methods depends on whether the directional feedback
30 information is available or not. Scenarios without directional feedback are confined to black-
31 box optimization methods such as evolutionary search [5], Bayesian optimization [6], or policy
32 gradient [9]. However, when the directional feedback is available, one can choose the much more
33 efficient gradient-based optimization method, such as stochastic gradient descent or exact line
34 search [1]. This insight motivates our use of directional feedback in the realm of LLM-driven
35 optimization.

36 As we show, the presence or absence of directional feedback and the possibility to access them is
37 crucial for LLM-based optimization. For a systematic study of factors that make an optimization

38 process challenging, we choose one of the difficult tasks proposed in the work on OPRO (listed
 39 as failure cases in Appendix A of OPRO [11]) – navigating a bumpy loss function landscape. We
 40 discover that an LLM-based optimizer agent’s performance varies with the information the feedback
 41 carries, and, given proper feedback, LLMs can strategically improve over past outputs, which makes
 42 this previously unsolvable task solvable. In addition, we demonstrate that using LLMs to “synthesize”
 43 feedback from a history of observations and prompts can help optimization too.

44 We also explore LLMs’ optimization potential in a completely different setting. Given the importance
 45 of feedback type in the LLM-based optimization process and the lack of benchmarks that generate
 46 verbal feedback automatically, we create a synthetic poem writing environment, where one can
 47 programmatically create feedback for the LLMs. The poem environment is a family of tasks where
 48 an LLM is asked to write a poem. A distinguishing feature of this benchmark is that the poems must
 49 satisfy some constraints, such as the number of syllables per line. By leveraging and synthesizing
 50 feedback, we show that an LLM can sequentially optimize a poem-generation prompt to yield a
 51 high success rate of producing constraint-satisfying poems. Our results highlight the importance
 52 of understanding and studying the role of feedback in the broader LLM-based text optimization
 53 landscape.

54 2 Preliminaries: Prompt Optimization for LLM-based Agent

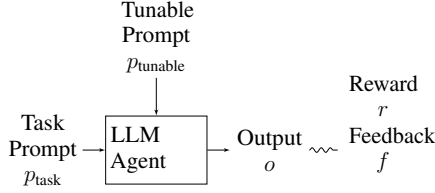
55 An LLM-based agent’s behavior is modulated through the prompts used as inputs to the LLM. We
 56 describe the interactive decision-making problem encountered by an LLM-based agent, and how
 57 prompt optimization through interactive feedback can improve the agent over time. In the following,
 58 uppercase letters, e.g. X , denote random variables or sets. Lowercase letters denote realizations of
 59 the random variables or set elements, e.g. “ $X = x$ ” states that a r.v. X takes on value x . Greek
 60 letters, e.g., ξ , denote parameters indexing probability distributions.

61 Consider an agent encountering a complex task such as generating a poem with logi-
 62 cal constraints. The task is communicated to the agent via a text prompt $p_{\text{task}} =$
 63 “Generate a poem with a rhyming scheme of abcba”. The LLM produces output text $o_1 \sim \text{Pr}_{\tau}(O \mid$
 64 $p_{\text{task}}, p_{\text{tunable}})$ (e.g., a poem, or plans, or executable code, or other texts as prompted), where τ cap-
 65 tures LLM hyper-parameters like sampling temperature, and p_{tunable} contains any orchestrated text
 66 inputs from other modules surrounding the LLM. In the sequel, we will develop modules that will
 67 incorporate information gathered over time to update p_{tunable} . By analogy with tunable parameters of
 68 an ML model, we view p_{tunable} as the tunable “parameters” of an LLM-based agent.

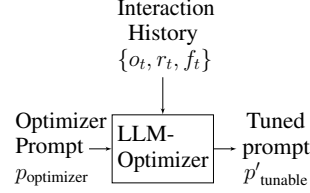
69 Based on the generated output o_1 , a scalar reward r_1 and optionally feedback f_1 are generated from
 70 the environment (e.g., human user response, or logs generated by executing code in a programming
 71 environment) and passed to the agent. For ease of notation, we assume $r \sim R$ and $f \sim F$, but we
 72 do not make specific assumptions of the underlying distributions. The reward can be a task success
 73 or failure boolean from the environment, or user-provided thumbs-up/down signal. This interaction
 74 process iterates $o_1 \rightsquigarrow \{r_1, f_1\}, \dots, o_t \rightsquigarrow \{r_t, f_t\}$, until the environment terminates the interaction
 75 session. Figure 1a illustrates the interaction process; for example, in Minecraft Voyager [10] a prompt
 76 $p_{\text{task}} =$ “Build a house” is translated into a code-generation request using internal orchestration
 77 that prepends a specific p_{tunable} ¹. The produced code o_t is executed in the Minecraft environment to
 78 generate error/debug/return messages f_t as well as task completion flag r_t that are returned to Voyager
 79 to refine the code o_{t+1} in subsequent iterations. The interaction session ends when the user prompting
 80 the Voyager agent terminates it. Note that p_{task} can be interactively updated within a session (e.g.,
 81 user providing additional hints or rephrasing the task), and we only assume that the rewards and
 82 feedbacks observed are consistent with the task that the agent is prompted to solve.

83 The LLM-Optimizer is a specific instance of an LLM-based agent. It can be used to improve another
 84 LLM-based agent using collected experience so that the generated outputs have higher expected
 85 reward $\mathbb{E}_o[r \mid o, p_{\text{task}}]$. The LLM-Optimizer takes a collection of Output-Reward-Feedback (o, r, f)
 86 tuples via its tunable prompt (see Figure 1b), and is tasked with generating a prompt p'_{tunable} for the
 87 LLM-based agent.

¹In Voyager, these prompts are hand-engineered and not automatically tuned.



(a) Schematic of LLM-based agent. $p_{tunable}$ can be updated from feedback and/or previous experiences via our sequential optimization.



(b) LLM-Optimizer is a specific LLM-based agent that incorporates previous experiences into the tunable prompt $p_{tunable}$ of the agent.

88 3 Optimizing LLM-based Agents

89 Define an LLM agent as $\pi : P_{task} \times P_{tunable} \rightarrow O$. The distribution O is defined by $p_{tunable}$ alone,
 90 which we can regard as the parameter of the LLM agent. The optimization problem we need
 91 to solve is to find $p_{tunable}^* := \arg \max_{p_{tunable}} \mathbb{E}_o [r \mid \pi(p_{task}, p_{tunable})]$. We can define an optimizer
 92 $g : P_{tunable} \times F \times R \rightarrow P_{tunable}$, where the goal is to find $p_{tunable}^*$ through a limited number of times
 93 that π attempts the task. An optimal optimizer g^* can find $p_{tunable}^*$ with the fewest amount of attempts.
 94 Different from the reasoning task setup, the output is defined as a distribution even for a single
 95 task p_{task} , and oftentimes, we do not know the distribution o^* that can obtain the highest expected
 96 reward.

97 3.1 Fundamentals of LLM Optimization

98 The most common approach for optimization is through an iterative solver that improves mono-
 99 tonically. However, in order to construct an iterative solver, the optimization problem needs to
 100 satisfy a few assumptions. To establish intuitions, we start with numerical optimization in a function
 101 approximation-based supervised learning setting. Given a hypothesis h and (x, y) , let $\tilde{y} = h_\theta(x)$.
 102 With a loss function $\ell : X \times Y \times \Theta \rightarrow \mathbb{R}$, we can define $L(\theta) = \mathbb{E}_{(x,y)} [\ell(\theta, x, y)]$. The goal is to
 103 find $\theta^* = \arg \min_\theta L(\theta)$. To achieve this goal, a valid optimization procedure proposes a new θ for
 104 k number of times. They usually consist of two steps:

105 **S1 Finding Valid Search Direction:** We need to find useful information, such as a descend direction
 106 $\Delta\theta^{(k)}$ that can help the update step. The usefulness of the information is tightly coupled with
 107 what the update step is.

108 **S2 Decide Update Rules:** We need to decide how to update θ . A typical update procedure is simply:
 109 $\theta^{(k+1)} = \theta^{(k)} + t^{(k)} \Delta\theta$, if $\Delta\theta$ is informative, where t_k is the step size.

110 If L is convex, then the criteria to determine whether $\Delta\theta$ is informative is quite simple: we can
 111 use the gradient of L . From convexity, we know that $\nabla L(\theta^{(k)})^T (L(\theta^{(k+1)}) - L(\theta^{(k)})) \geq 0$ implies
 112 $L(\theta^{(k+1)}) \geq L(\theta^{(k)})$. Then we can set the descend direction $\Delta\theta^{(k)}$ to satisfy $-\nabla L(\theta^{(k)})^T \Delta\theta^{(k)}$.
 113 A simple way to satisfy this criterion is let $\Delta\theta^{(k)} := -\nabla L(\theta^{(k)})$, which is the gradient descent
 114 method (GD). However, more complicated update rules can be used, such as backtracking line
 115 search [1]. It is also worth noting that we do not always need to satisfy **S2**. For example, in an
 116 evolutionary search algorithm, many candidates are proposed and the update rule is simply to keep
 117 the candidate with the best score.

118 Then we can contrast the setting with an LLM optimization problem. If we want to have an iterative
 119 descent algorithm to find the optimal prompt for an LLM agent, then we need to consider the
 120 following properties:

121 **S1 Search Direction:** We should obtain useful information, analogous to $\nabla L(\theta)$, to help inform the
 122 optimizer on how to update the parameter $p_{tunable}$.

123 **S2 Update Parameter:** Unlike the numerical case, where basic algebra can be applied to update
 124 parameters, it is unclear whether there is a predefined notion of $\Delta p_{tunable}^{(k)}$ in text space. This
 125 distance $\Delta p_{tunable}^{(k)}$ in the best case, can be assessed by human intuition over the semantics of the
 126 text, in the worst case, can be completely arbitrary.

127 In order to propose an optimization algorithm using LLM as an optimizer, we must make the following
 128 assumptions:

129 **A1 Permissible Search Direction:** There exists useful information, which we describe as feedback,
 130 f , for an LLM optimizer g such that g can propose a p_{tunable}^{k+1} where $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^{k+1})] \geq$
 131 $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^k)]$.

132 **A2 Valid Update:** The LLM optimizer g can modify p_{tunable} based on f , where direction of change:
 133 $\Delta p_{\text{tunable}}$ is determined by the information contained in f (i.e., not a random text edit).

134 In the next few sections, we describe a few possible settings where these assumptions can be satisfied
 135 or need not be. We assume **A2** is always satisfied. The first setting we discuss is that it is possible
 136 that LLM itself acts as a black-box optimizer. For example, it might obtain a valid search direction
 137 by implicitly computing finite differences between inputs and outputs in text space.

138 **LLM Might Implicitly Perform Newton’s Method** Similar to Newton’s Method for using finite
 139 difference to approximate gradient, we can hope that LLM can implicitly compute the following
 140 function:

$$\nabla R = \lim_{\Delta p_{\text{tunable}}^{(k)} \rightarrow 0} \frac{\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^k)] - \mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^{k+1})]}{\Delta p_{\text{tunable}}^{(k)}}$$

141 If we think this is possible, then the input to the LLM can be tuples of $(p_{\text{tunable}}^1, r_1), \dots, (p_{\text{tunable}}^k, r_k)$.
 142 Although it is unclear if this is truly the case, this shows that the optimizer needs to retain a history of
 143 how past prompts p have changed the reward r .

144 **Feedback Can Be Directional** The other possibility, not relying on the black box “magic” of the
 145 LLM’s internal process, is to hope that somehow a permissible search direction f is given to us from
 146 an external source. Humans give directional feedback quite often: “This coffee is too hot for me.” or
 147 “Can you lower the room temperature?” The first feedback implicitly asks the agent to make a cooler
 148 coffee (but not saying exactly how cool it should be). The second feedback asks the agent to turn down
 149 the room temperature (but without specifying which temperature to set). Imagine if the agent’s action
 150 (output space O) for both cases is to write API calls to set temperature; then we know immediately
 151 what ΔO should be – keep everything the same, but enter a lower temperature value. After the
 152 adjustment, a user might say: “This coffee is now too cold for me.” or “I’m freezing!” Making this
 153 kind of feedback very similar to the gradient information we get from numerical optimization. This
 154 suggests that, in some cases, incorporating feedback (or somehow obtaining directional feedback)
 155 could be helpful for the optimization procedure. We should provide LLM optimizer with examples of
 156 $(p_{\text{tunable}}^1, f_1, r_1), (p_{\text{tunable}}^2, f_2, r_2), \dots, (p_{\text{tunable}}^k, f_k, r_k)$ instead.

157 **Non-directional Feedback** There is another type of feedback we can consider. This type of
 158 feedback contains useful information but is not directional because they do not directly inform us how
 159 to change the input O . For example, feedback like “I can’t drink this coffee because the temperature
 160 is not quite right.” This feedback clearly states that the attribute of “temperature” is important to the
 161 user and is not satisfactory. However, it does not tell us whether we should make a coffee that’s hotter
 162 or colder. Coffee can have many attributes, such as “temperature”, “acidity”, “roast”, “sweetness”, or
 163 “cream-level”. This feedback is more useful than a scalar reward because it *explains* attributes that
 164 affect R and allows us to focus more on a single dimension of many attributes.

165 **Reward as Feedback / No Feedback** Unlike directional and non-directional feedback, which
 166 usually contains information about how to change p_{tunable} , score-based feedback only gives back a
 167 numerical value indicating how well p_{tunable} performs. In our setup, this means LLM-Optimizer only
 168 observes reward r without f . This is often referred to as the 0th-order feedback.

169 3.2 Sequential Prompt Optimization

170 Inspired by the descent method in numerical optimization, we propose an algorithm that aims to
 171 satisfy the requirement of descent methods such that we can reach the extremum. We define the
 172 following optimization loop with an LLM-based agent π . Agent π with an initial tunable prompt
 173 p_{tunable} takes a task description p_{task} from the environment and samples an output o_1 . The environment
 174 returns a reward r_1 and feedback f_1 . An LLM-optimizer stores $(o_1, r_1, f_1, p_{\text{tunable}})$ in a history buffer
 175 \mathbb{H} . When the buffer becomes large, we subsample H from \mathbb{H} . The LLM-optimizer proposes a
 176 new tunable parameter p'_{tunable} . We make an explicit decision on whether to replace p_{tunable} with
 177 p'_{tunable} based on the reward evaluated on the distribution o and o' . We describe the full procedure
 178 in Algorithm 1. We now describe the implementation choices for each component of our iterative
 179 solver.

180 **Policy** Policy is an LLM that takes in a tunable instruction prompt p_{tunable} , description of task p_{task}
 181 and produces an output o . It does not see a history of interaction it did with the environment. This
 182 design decision prohibits the policy from controlling its own prompts. It is different from some other
 183 existing work. For example, Voyager [10] would allow the policy to see all of its interaction histories
 184 (and errors they make). React [12] also allows the policy to see the full error trace. Allowing the
 185 policy to see its past errors is a specific design choice on p_{tunable} . We hope to let the optimizer decide
 186 what p_{tunable} should be without injecting human prior.

187 **Prompt Proposal** We define the prompt proposal module as $\Delta : P_{\text{task}} \times P_{\text{tunable}} \times F \times R \rightarrow P_{\text{tunable}}$.
 188 This module looks at the task description, past prompts, and the feedback each prompt receives and
 189 proposes a new prompt. If the environment provides directional or non-directional feedback, this
 190 module should take in F as well. Even though past prompts were included, this module is allowed to
 191 produce completely new prompts.

Algorithm 1: Sequential Prompt Optimization

Input: Given s, R , an LLM-based agent π , an
 LLM-based prompt proposal module Δ , p_{tunable}^0 ,
 p_{task} , and K iterations.

Output: p_{tunable}^K

$\mathbb{H} = \emptyset$

for $k \leftarrow 0 \dots K$ **do**

$o_k, r_k \sim \pi(p_{\text{task}}, p_{\text{tunable}}^{(k)}), R$
 $f_k \sim F$ or $\hat{F}(p_{\text{task}}, o_k, r_k)$
 $H = \text{Sample}(\mathbb{H})$
 $p^{(k+1)} = \Delta(H, \{p_{\text{task}}, p_{\text{tunable}}^k, f_k, r_k\})$
if $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^{k+1})] \geq$
 $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^k)]$ **then**
 | $p_{\text{tunable}}^{k+1} = p_{\text{tunable}}^k$
end
 $\mathbb{H} = \mathbb{H} \cup \{o_k, r_k, f_k, p_{\text{tunable}}^k\}$

end

return p_{tunable}^K

193 **Feedback Synthesizer** If numerical feedback is the only type of feedback given, we can design
 194 a feedback synthesizing module $\hat{F} : P_{\text{task}} \times O \times R \rightarrow F$. It takes in $(o_1, r_1), \dots, (o_k, r_k)$ and
 195 produces feedback to the prompt proposal module. We designed this module to ask the question,
 196 “How should the input be changed to have a greater effect on the objective/output?” We note that we
 197 more specifically prompt the LLM to think about the difference in the input space that would impact
 198 the difference in output space, different from APO, where they ask the LLM to “give reasons why the
 199 prompt could have gotten these examples wrong.”

200 **Prompt Selector** In order to satisfy the descent method assumption **A1** (permissible search di-
 201 rection), we need to guarantee that p_{tunable}^{k+1} is an improvement over p_{tunable}^k . This can be achieved by
 202 setting a selection criterion that requires $p^{(k+1)}$ to get a higher reward than $p^{(k)}$. A simple crite-
 203 rion is to improve the average performance over the distribution of O : $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^{k+1})] \geq$
 204 $\mathbb{E}_o [r \mid \pi(p_{\text{task}}, p_{\text{tunable}}^k)]$.

205 4 Experiments

206 4.1 Numerical Optimization

207 We test if LLM is possible to do optimization and what are the necessary ingredients for it to find the
 208 optimal solution in an optimization problem. We set up this task because in prompt optimization,
 209 it is often hard to know what the optimal prompt or a good search direction is. With a numerical
 210 optimization problem, both are well-defined. We use a set of classic optimization problems² that
 211 require LLMs to find x , a 2-dimensional vector.

²<https://www.sfu.ca/~ssurjano/optimization.html>

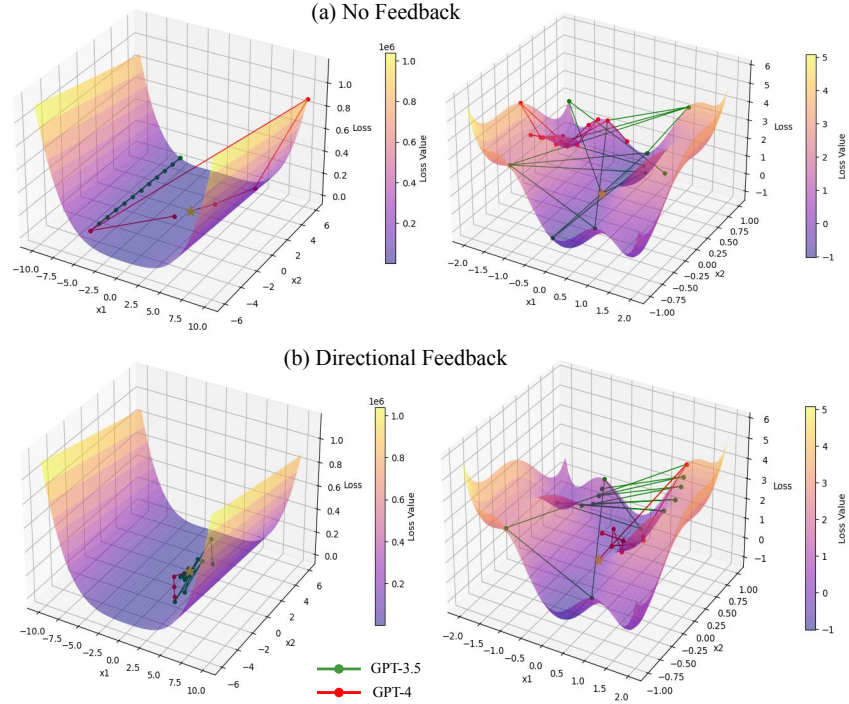


Figure 2: We visualize the optimization trajectory path made by the Optimizer Agent with GPT-3.5 and GPT-4. The loss landscape on the left is the Rosenbrock Function, and on the right is the Six-Hump Camel Function.

- 212 1. Task: Given a task description p_{task} and a function J (which is hidden from the LLM), we sample
- 213 a random starting point $(x_0, J(x_0))$, $x_0 \sim X$. An LLM is asked to produce x to minimize J .
- 214 2. Optimizable variable: X . The LLM is asked to output x directly. Here, p_{tunable} is the same as x .
- 215 3. Output process O : the output module takes x and directly outputs x , an identity function.
- 216 4. Reward R : $R(x) = -J(x)$.
- 217 5. Feedback F :
 - 218 • Directional Feedback: $\nabla R(x) = \frac{dR}{dx}$, the first-order derivative of the output.
 - 219 • Non-directional Feedback: We compare the partial derivatives $\frac{\partial R}{\partial x_1}$ and $\frac{\partial R}{\partial x_2}$, and tell the LLM
 - 220 which dimension of x should be changed to accomplish the task (but without telling LLM in
 - 221 which direction to change).

222 This experiment does not use the full setup of Algorithm 1. We only test our LLM-based optimizer Δ

223 and our feedback synthesizer \hat{F} . We choose four functions: Booth, McCormick, Rosenbrock, and

224 Six-Hump-Camel Function. They were chosen because the optimal x that minimizes these functions

225 is not $[0, 0]$. In our initial experiments, LLM is quick to guess $[0, 0]$, which trivializes the optimization

226 problem.

227 We define simple regret as $\text{Reg}(\Delta) = |J(x_T) - J(x^*)|$, where J is the function we try to minimize

228 and T is the number of optimization steps we allow the optimizer Δ to take. We define cumulative

229 regret as $\text{CuReg}(\Delta) = \sum_{t=1}^T |J(x_t) - J(x^*)|$. Intuitively, simple regret corresponds to how close is

230 an optimizer’s final answer x_T to the correct answer x^* . Cumulative regret describes how “efficient”

231 is the optimizer at finding the x^* . We compare three models: Δ with **GPT-3.5**, **GPT-4**, and a

232 stochastic gradient descent algorithm (**SGD**) with a small yet fixed learning rate. In the reported

233 results, we run 10 trials and allow Δ to take at most 10 optimization steps.

234 **RQ1** Can an LLM Implicitly Perform Newton’s Method, given $(x_1, J(x_1)), \dots, (x_k, J(x_k))$?

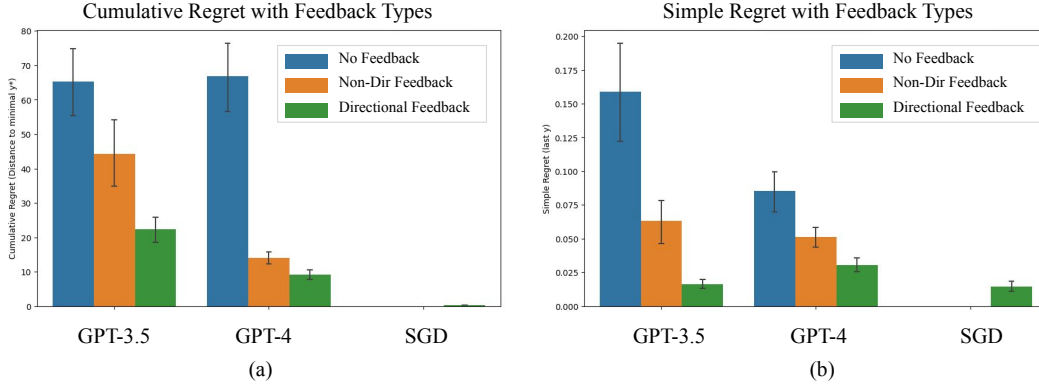


Figure 3: We plot out the average Cumulative Regret and Simple Regret of each condition over 10 trials. Each algorithm is allowed to take 10 steps. We tuned the SGD learning rate slightly to ensure it is not too large or too small. The result is aggregated over 4 loss functions.

235 From Figure 2, we can see that LLM, as an optimizer, has a rough sense of direction, given a history
 236 of past explorations. In Figure 2 (a), we note that in both loss landscapes, although GPT-3.5 often
 237 fails to find the minimal point without feedback (green lines), GPT-4 is able to understand the past
 238 history and make new proposals of x that incrementally minimizes $J(x)$. This suggests that even
 239 though there is no explicit gradient computation, LLM can be asked to “improve” based on a history
 240 of observations.

241 **RQ2** Does directional Feedback help the optimization process? Do other types of feedback help as
 242 much?

243 We designed the prompt space for the LLM-based optimizer Δ to insert feedback text right after the
 244 observation text and with an additional wording that reads, “You should incorporate the suggestion.”
 245 Besides this change, the optimizer agent prompt stays the same between no feedback and with
 246 feedback conditions. The full prompt is available in the Appendix.

247 From Figure 2 (b) and Figure 3, we can see that both GPT-3.5 and GPT-4 are able to take advantage
 248 of the additional feedback information and improve their search direction. Feedback can help both a
 249 weaker model (GPT-3.5) and a strong model (GPT-4). A stronger model can improve more, even
 250 if the feedback has less information (see the comparison between Non-directional Feedback and
 251 Directional Feedback in Section 3.1).

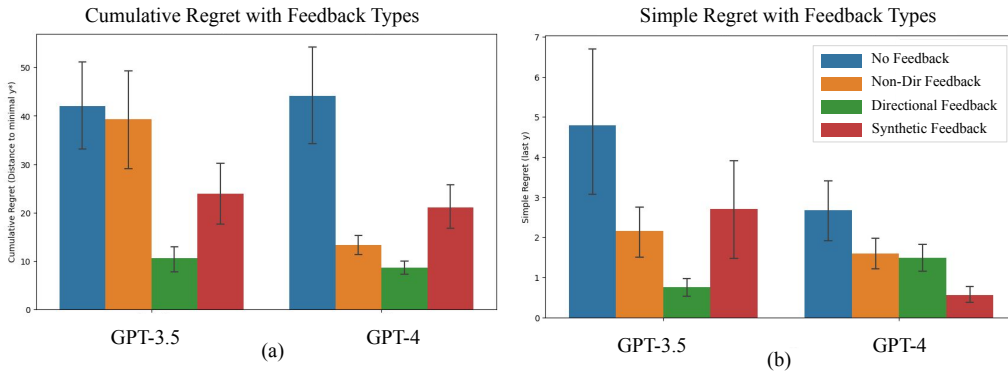


Figure 4: We plot out the average Cumulative Regret and Simple Regret of each condition over 10 trials and compare different feedback types. **Synthetic Feedback** is generated by the same LLM as the optimizer.

252 Although loss minimization is a challenging task for LLMs, with some amount of feedback, LLMs
 253 are able to find a final x that is similar to a classic optimization algorithm like SGD (see Figure 3b)
 254 – the simple regret is similar. It is worth noting that GPT-4’s final proposed x is not as close to the
 255 optimal as GPT-3.5. This is potentially because both models decide their own step size, and we are
 256 limiting the optimization horizon to 10 steps.

257 **RQ3** If directional feedback is missing, can we replace it with an LLM module to enhance whatever
 258 feedback is available?

259 Oftentimes, direct and useful feedback might be missing from the environment. In this experiment,
 260 we design a feedback synthesizer module (described in Sec 3.2) that can take the output from the
 261 model and the reward and try to provide feedback that can improve the next output. Different from
 262 methods such as self-reflection, self-criticism, or thinking step-by-step, the feedback synthesizer asks
 263 questions similar to “How should I change about x that will result in a larger change in y ?”, where
 264 self-reflection usually asks the model to reflect upon past “mistakes” on what they did wrong.

265 In Figure 4, we show that we can synthesize feedback from a history of past outputs and rewards
 266 that is able to guide the optimizer LLM to find a better solution. Synthesized feedback is not as
 267 informative as directional feedback that comes from the environment, but it can easily outperform
 268 settings where no feedback is given.

269 4.2 Poem Generation

270 Now, we have validated the importance of feedback. We want to validate our optimization setup on a
 271 more challenging domain, where now we have to optimize over a prompt that controls how another
 272 LLM-based agent produces output. An easy constrained optimization problem to set up is poem
 273 generation. A formal poem is a writing assignment that requires the creation of a poem to satisfy
 274 some requirements regarding its form. For example, Haiku is a type of formal poem that asks for
 275 three lines that form a 5-7-5 syllable pattern. This is a challenging task for both GPT-3.5 and GPT-4,
 276 but easy for us to verify whether the generated poem has satisfied the constraint.

- 277 1. Task: Generate a poem with a given constraint sampled from a set of constraints.
- 278 2. Optimizable variable: P_{tunable} : This is the prompt p_{tunable} for the LLM-based agent that we want to
 279 update and optimize.
- 280 3. Output process O : the LLM agent takes the prompt p_{tunable} and follows its suggestion and a task
 281 description p_{task} to produce a poem o .
- 282 4. Reward R : The fraction of lines in the generated poem that satisfy the constraint described by
 283 p_{task} . $r \in [0, 1]$.
- 284 5. Feedback F :
 - 285 • Directional Feedback: We print out the number of syllables in the current line and whether
 286 LLM needs to increase or decrease the number of syllables in that line.
 - 287 • Non-directional Feedback: We print out how many lines violate the poem writing constraints.

288 In the following experiment, we use the full setup of Algorithm 1. We allow each agent to take 10
 289 optimization steps. We name our agent **OptAgent**. It produces an instruction that will be sent to the

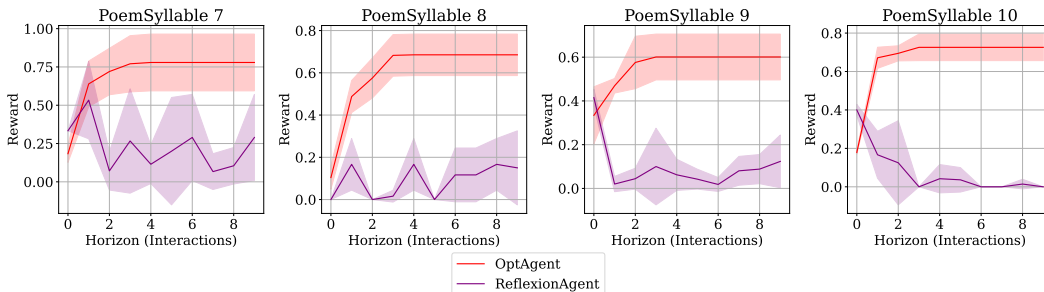


Figure 5: We show the reward for each policy after each round of interaction with the environment. OptAgent (our algorithm) is in red.

290 poem generation agent to produce a poem. The poem-generation agent will not see the history of
291 mistakes or any other information. We additionally evaluate **Reflexion agent** [8]. We set up four
292 tasks: generating poems that contain 7, 8, 9, or 10 syllables for each line.

293 We show that in Figure 5, we can reliably select prompts that improve the policy performance for
294 each task. The prompt selection step in our optimization algorithm ensures that the new prompt
295 will improve the performance. Otherwise, it will reject the updated prompt and keep the previous
296 prompt. This differs from the Reflexion Agent in that they are not guaranteed to improve in the next
297 interaction.

298 5 Conclusion

299 This paper argues that LLMs can successfully optimize a wide range of entities ranging from from
300 mathematical functions to prompts for textual tasks if provided with directional feedback. We
301 empirically show on challenging numerical optimization scenarios and constrained text generation
302 tasks that utilizing either environment-provided or synthetic feedback is a crucial piece in LLM-based
303 optimization. We emphasize that this is an early work on general LLM-based optimizers. LLMs'
304 potential in this role is still waiting to be realized with improved methods for directional feedback
305 generation.

306 References

- 307 [1] Boyd, S. P. and Vandenberghe, L. (2004). *Convex optimization*. Cambridge university press.
- 308 [2] Gur, I., Furuta, H., Huang, A., Safdari, M., Matsuo, Y., Eck, D., and Faust, A. (2023). A
309 real-world webagent with planning, long context understanding, and program synthesis. *arXiv*
310 *preprint arXiv:2307.12856*.
- 311 [3] Ichter, B., Brohan, A., Chebotar, Y., Finn, C., Hausman, K., Herzog, A., Ho, D., Ibarz, J., Irpan,
312 A., Jang, E., Julian, R., Kalashnikov, D., Levine, S., Lu, Y., Parada, C., Rao, K., Sermanet, P.,
313 Toshev, A. T., Vanhoucke, V., Xia, F., Xiao, T., Xu, P., Yan, M., Brown, N., Ahn, M., Cortes, O.,
314 Sievers, N., Tan, C., Xu, S., Reyes, D., Rettinghouse, J., Quiambao, J., Pastor, P., Luu, L., Lee,
315 K.-H., Kuang, Y., Jesmonth, S., Joshi, N. J., Jeffrey, K., Ruano, R. J., Hsu, J., Gopalakrishnan,
316 K., David, B., Zeng, A., and Fu, C. K. (2023). Do as i can, not as i say: Grounding language in
317 robotic affordances. In *CORL*.
- 318 [4] Kojima, T., Gu, S. S., Reid, M., Matsuo, Y., and Iwasawa, Y. (2022). Large language models are
319 zero-shot reasoners. In *NeurIPS*.
- 320 [5] Mitchell, M. (1998). *An introduction to genetic algorithms*. MIT press.
- 321 [6] Mockus, J. (1998). The application of bayesian methods for seeking the extremum. *Towards*
322 *global optimization*, 2:117.
- 323 [7] Pryzant, R., Iter, D., Li, J., Lee, Y. T., Zhu, C., and Zeng, M. (2023). Automatic prompt
324 optimization with "gradient descent" and beam search. *arXiv preprint arXiv:2305.03495*.
- 325 [8] Shinn, N., Cassano, F., Labash, B., Gopinath, A., Narasimhan, K., and Yao, S. (2023). Reflexion:
326 Language agents with verbal reinforcement learning. *arXiv preprint arXiv:2303.11366*.
- 327 [9] Sutton, R. S., McAllester, D., Singh, S., and Mansour, Y. (1999). Policy gradient methods for
328 reinforcement learning with function approximation. In *NeurIPS*.
- 329 [10] Wang, G., Xie, Y., Jiang, Y., Mandlekar, A., Xiao, C., Zhu, Y., Fan, L., and Anandkumar, A.
330 (2023). Voyager: An open-ended embodied agent with large language models. *arXiv preprint*
331 *arXiv:2305.16291*.
- 332 [11] Yang, C., Wang, X., Lu, Y., Liu, H., Le, Q. V., Zhou, D., and Chen, X. (2023). Large language
333 models as optimizers. *arXiv preprint arXiv:2309.03409*.
- 334 [12] Yao, S., Zhao, J., Yu, D., Du, N., Shafran, I., Narasimhan, K. R., and Cao, Y. (2023). React:
335 Synergizing reasoning and acting in language models. In *ICLR*.

336 **A Appendix**

337 **A.1 Loss Optimizing Experiment Details**

338 We designed the prompt for two agents. The prompt is written in a Handlebar syntax, where “`{{}}`”
339 indicate variables to be replaced. A brief guide on this syntax is available here³.

340 **For the LLM-Optimizer:**

```
341 {{#system~}}
342 You are trying to minimize the output (y) of a function by choosing input (x).
343 {{~/system~}}
344
345 {{#user~}}
346 {{task_description}}
347
348 This is what you have previously chosen for x and what the ys were:
349 {{observation}}
350
351 {{feedback}}
352
353 You should incorporate the suggestion to output the next x.
354 Please output the next x that will make this function output the smallest y.
355 You cannot repeat the same x, doing so will result in a penalty.
356
357 Format: x = [x1, x2]
358 Output:
359 {{~/user~}}
```

360 **For the feedback synthesizing LLM:**

```
361 {{#system~}}
362 You are trying to minimize the output (y) of a function by choosing input (x).
363 {{~/system~}}
364
365 {{#user~}}
366 You are trying to minimize the output (y) of a function by choosing input (x).
367 You get to observe y once you choose the value of x, where x is a 2-dimensional vector.
368 This means x = [x1, x2], where x1 and x2 are real numbers.
369 The goal is to choose x such that y is as small as possible.
370
371 Here is a list of x and how it affects y:
372 {{#each history}}
373 {{this.action}}
374 {{this.observation}}
375 =====
376 {{~/each}}
377
378 For x = [x1, x2]
379 What are the suggestions you can give to the user to make y smaller?
380 For example, here are some of the things you can suggest:
381 - Changing x1 seems to have a bigger effect on y than changing x2.
382 - Make a larger change on x2
383 - Increase x1 by 1.2
384 - Decrease x2 by 0.5
385 - Try to increase x1 and decrease x2 at the same time
386 Or any other kind of suggestion. Do not make a suggestion that's the form of a question.
387 You should only make a one-sentence suggestion that's brief and short.
388
389 Suggestion:
390 {{~/user~}}
```

391 **A.2 Poem Experiment Details**

392 **For the LLM-agent that generates the poem, we use the following prompt:**

³<https://github.com/guidance-ai/guidance>

```

393 {{#system~}}
394 You are a student and your teacher gives you an assignment to write a poem.
395 {{~/system}}
396
397 {{#user~}}
398 The assignment is:
399 {{assignment}}
400
401 {{#if exists_intrusction}}
402 In addition, here are some helpful advice and guidance:
403 {{instruction}}
404 {{/if}}
405 {{~/user}}

```

406 **For the feedback synthesizer module, we use this prompt:**

```

{{#system~}}
You are a helpful assistant who aims to provide feedback to a student
↳ who's writing a poem
according to some instructions.
It is important to let the student know if they did satisfy the
↳ instruction or not and why.
{{~/system}}

{{#user~}}
This is the history of past generated poems and how well they did with
↳ respect to instructions.

{{#each history}}
Instruction: {{this.observation}}

Poem:
{{this.action}}

Feedback from the teacher:
{{this.feedback}}
-----
{{~/each}}
{{~/user}}

{{#user~}}
Now, the student writes a new poem.

New instruction: {{observation}}

Poem:
{{action}}

What changes can you make to the poem to help it conform to the
↳ instructions?
{{~/user}}

{{#assistant~}}
{{gen 'exp_feedback' temperature=0.7}}
{{~/assistant}}

```

407 **For the LLM-based optimizer, we use this prompt:**

```

{{#system~}}
You are a helpful assistant that wants to come up with instructions to a
↳ student to help them write a poem that is satisfactory to a teacher's
↳ assignment.
The student's poem needs to satisfy the requirement of this assignment.
{{~/system}}

{{#user~}}

```

This is the history of how you have been helping this student and whether
↳ your instructions have succeeded.
Teacher's feedback is the most important feedback, because the student
↳ needs to meet the teacher's criteria.
However, another student's feedback can provide helpful information too.

{{#each history}}
The Assignment: "{{this.assignment}}"

Your Instruction:
{{this.prompt}}

Student's Poem:
{{this.action}}

Teacher's Feedback:
{{this.feedback}}

Feedback from another student:
{{this.exp_feedback}}

{{~/each}}
{{~/user}}

{{#user~}}

Your previous instruction didn't work -- the students didn't write a poem
↳ that satisfied the teacher's criteria.
Based on your interaction with the students, can you come up with better
↳ instructions that can help this student write a poem that matches the
↳ teacher's criteria?
Keep in mind that telling the student what to do step-by-step might be
↳ very helpful!
However, you need to be brief and to the point.
{{~/user}}