
1 CAPO - Pseudo Algorithm

Algorithm 1 CAPO: Cost-Aware Prompt Optimization

Require: Dataset $\mathcal{D} = \{(X_i, y_i)\}_{i=1}^n$, Meta-LLM $\Phi(x)$, Downstream LLM $\phi(x)$, Initial instructions $\Lambda = \{\lambda_1, \dots, \lambda_p\}$, Population size p , Block size b , Number of iterations n , Number of crossovers per iteration c

```
1: for  $\lambda \in \Lambda$  do
2:   num_shots  $\sim U(\text{lower}, \text{upper})$ 
3:    $\xi \leftarrow \text{sample}(X, \text{num\_shots})$  ▷ Sample few shot examples
4:    $\theta \leftarrow \phi(\lambda || \xi)$  ▷ Generate few shots with reasoning
5:    $\pi \leftarrow (\lambda, \theta)$ 
6:    $\Pi \leftarrow \Pi.\text{append}(\pi)$ 
7: end for
8: Divide dataset  $\mathcal{D}$  into blocks  $\mathcal{B} = \{B_1, \dots, B_k\}$  where  $|B_i| = b$ 
9: for  $i = 1$  to  $n$  do
10:   $\Pi_{\text{off}} \leftarrow \text{cross\_over}(\Pi, c)$ 
11:   $\Pi_{\text{off}} \leftarrow \text{mutate}(\Pi_{\text{off}})$ 
12:   $\Pi \leftarrow \text{do\_racing}(\Pi \cup \Pi_{\text{off}}, k = p)$ 
13: end for
14: best_prompt  $\leftarrow \text{do\_racing}(\Pi, k = 1)$ 
15: return best_prompt
```

Is reasoning for few-shots always necessary? Maybe this is not important for simple tasks and costs a lot of tokens. (Idea: some few shot examples can be included without reasoning)

Algorithm 2 cross_over

Require: Population Π , Meta-LLM $\Phi(x)$, Cross-Over-Meta-Prompt λ_C , Number of crossovers c

```
1:  $\Pi_{\text{off}} \leftarrow []$ 
2: for  $j = 1$  to  $c$  do
3:    $P_1 \leftarrow \text{sample}(\Pi, 1)$   $\triangleright P_1 = (\lambda_1, \theta_1)$ 
4:    $P_2 \leftarrow \text{sample}(\Pi, 1)$   $\triangleright P_2 = (\lambda_2, \theta_2)$ 
5:    $\lambda_{\text{off}} \leftarrow \Phi(\lambda_C || \lambda_1 || \lambda_2)$   $\triangleright$  Let Meta-LLM cross over the prompts
6:    $\theta_{\text{off},j} \leftarrow \text{sample}(\theta_1 \cup \theta_2, \left\lfloor \frac{|\theta_1| + |\theta_2|}{2} \right\rfloor)$   $\triangleright$  Sample from all few-shot examples
7:    $\pi_{\text{off},j} \leftarrow (\lambda_{\text{off},j}, \theta_{\text{off},j})$ 
8:    $\Pi_{\text{off}} \leftarrow \Pi_{\text{off}}.\text{append}(\pi_{\text{off},j})$ 
9: end for
10: return  $\Pi_{\text{off}}$ 
```

Currently we do a random parent selection. In EvoPrompt they do a roulette wheel selection based on the fitness scores. This would require us to find a way of already having scores here.

We have to clarify where the new few-shot examples are coming from.

Extra Split:

- + no data leakage (fair, comparable assessment of prompts)
- constrained pool of few-shot examples (how to get this pool?) - potentially smaller dev Split

From Train Split:

- data leakage (prompts that contain eval data point already as few-shot examples which are already confirmed as correct have advantages)
- + we can use the full train set

Algorithm 3 mutate

Require: Population of offsprings Π_{off} , Meta-LLM $\Phi(x)$, Mutation-Meta-Prompt λ_{M} ,
Dataset samples X

```
1: for  $\pi_{\text{off}} \in \Pi_{\text{off}}$  do
2:    $\lambda_{\text{off}} \leftarrow \Phi(\lambda_{\text{M}} || \lambda_{\text{off}})$ 
3:    $\text{num\_shots} \sim \text{U}(\text{lower}, \text{upper})$ 
4:    $\text{num\_new\_shots} \sim \text{U}(\text{lower}, \text{num\_shots})$ 
5:    $\xi \leftarrow \text{sample}(X, \text{num\_new\_shots})$   $\triangleright$  Sample new few shot examples
6:    $\theta_{\text{new}} \leftarrow \phi(\lambda_{\text{off}} || \xi)$   $\triangleright$  Generate few shots with reasoning
7:    $\theta_{\text{old}} \leftarrow \text{sample}(\theta_{\text{off}}, \text{num\_shots} - \text{num\_new\_shots})$ 
8:    $\theta \leftarrow \theta_{\text{old}} \cup \theta_{\text{new}}$ 
9:    $\theta \leftarrow \text{shuffle}(\theta)$ 
10:   $\pi_{\text{off}} \leftarrow (\lambda_{\text{off}}, \theta)$ 
11: end for
12: return  $\Pi_{\text{off}}$ 
```

Algorithm 4 do_racing

Require: Prompts Π , Top-k k , blocks \mathcal{B} , Downstream LLM $\phi(x)$, Max number of evaluated blocks max_n_blocks_eval

```
1:  $i \leftarrow 0$ 
2:  $\text{scores} \leftarrow [0] * \text{len}(\Pi)$ 
3:  $\text{shuffle}(\mathcal{B})$  ▷ Whether to shuffle is left as a HP
4: while  $\text{len}(\Pi) > k \wedge i < \text{max\_n\_blocks\_eval}$  do
5:    $i \leftarrow i + 1$ 
6:    $\text{scores} \leftarrow \frac{1}{i} (\text{evaluate}(\Pi, B_i) + (i - 1) * \text{scores})$  ▷ Already evaluated blocks are cached
7:    $\Pi \leftarrow \text{racing\_elimination}(\Pi, \text{scores}, \alpha, k)$ 
8: end while
9: if  $\text{len}(\Pi) > k$  then
10:    $\Pi \leftarrow \text{top\_k}(\Pi)$  ▷ Make sure to return only k prompts
11: end if
12: return  $\Pi$ 
```

Algorithm 5 racing_elimination

Require: Survivors Π , scores S , confidence level α , top-k k

```
1:  $c_\alpha \leftarrow \text{getCriticalValue}(\alpha)$ 
2: for  $\pi_i \in \Pi$  do
3:    $\text{n\_sig\_better} \leftarrow \sum_{j \neq i} \mathbf{1}_{[\text{getTestStatistic}(s_j, s_i) > c_\alpha]}$ 
4:   if  $\text{n\_sig\_better} \geq k$  then
5:      $\Pi \leftarrow \Pi \setminus \{\pi_i\}$  ▷ Eliminate  $\pi_i$ 
6:   end if
7: end for
8: return  $\Pi$ 
```

Potential test statistics:

- Friedmann test (with post-hoc tests) like mentioned in irace paper
- Two-sided t-test like mentioned in irace paper
- Hoeffding races: $\varepsilon = \sqrt{\frac{B^2 \log(2/\delta)}{2n}}$

Further Ideas:

Multi-objective optimization (predictive performance and cost(e.g. prompt length)) -

how to consider both objectives in the test statistic? -> only then we are “cost aware”
(also reduces cost of optimization because of shorter prompt, fewer few-shot examples,
...)

Multi-fidelity can be e.g. done by first using a smaller model to determine a promising population and then continue with a larger model.