# GenDLN: Evolutionary Algorithm-Based Stacked LLM Framework for Joint Prompt Optimization

Anonymous ACL submission

#### Abstract

With Large Language Model (LLM)-based 002 applications becoming more common due to strong performance across many tasks, prompt optimization has emerged as a way to extract better solutions from frozen, often commercial LLMs that are not specifically adapted to a task. LLM-assisted prompt optimization meth-007 ods provide a promising alternative to manual/human prompt engineering, where LLM "reasoning" can be used to make them optimizing agents. However, the cost of using 011 LLMs for prompt optimization via commercial APIs remains high, especially for heuristic 013 methods like evolutionary algorithms (EAs), which need many iterations to converge, and 015 thus, tokens, API calls, and rate-limited network overhead. We propose GenDLN, an effi-017 cient genetic algorithm-based prompt pair optimization framework that leverages commercial API free tiers. Our approach allows teams with limited resources (NGOs, non-profits, academics...) to efficiently use commercial LLMs for EA-based prompt optimization. We conduct experiments on CLAUDETTE for legal terms of service classification and MRPC for paraphrase detection, performing in line with 027 selected prompt optimization baselines, at no cost. Our code is available in **<omitted>**.

#### 1 Introduction

037

041

LLMs (large language models) are increasingly replacing traditional classification and inference models due to their generality, ability to perform a wide range of tasks, and seemingly advanced "reasoning." As the use of LLMs for domain-specific tasks becomes more ubiquitous, prompt optimization emerges as an important area of research to improve the task-specific performance of LLMs, especially in complex domains like legal text analysis and interpretation (Hakimi Parizi et al., 2023; Lai et al., 2024). In recent years, several prompt design and optimization techniques have been proposed.



Figure 1: Running an individual through the DLN, where an individual is a prompt pair  $(p_1, p_2)$ ;  $LLM_1$  responds to  $p_1$ , and this response, along with  $p_2$ , is fed into  $LLM_2$  for classification. E.g.:  $p_1$ : "Interpret <ToS sentence i>" -  $p_2$ : "Based on the above interpretation, classify <ToS sentence i> as fair or unfair."

Some examples are edit-based instruction search GRIPS (Prasad et al., 2023) and reflection-based frameworks that incorporate LLM self-critique such as ProTeGi (Pryzant et al., 2023) and OPRO (Yang et al., 2024).

Deep Language Networks (DLNs) is a novel approach that stacks LLMs as computational units (Sordoni et al., 2023). Like other prompt optimization methods, the goal is to use frozen-weight LLMs for inference while refining input prompts for better results. However, they stack two LLMs, jointly optimizing two input prompts, where the output of the first LLM, along with the second prompt, are fed into the second LLM, as shown in Fig. 1. The prompts are treated as learnable parameters of the generative distribution, and the prompt pair is jointly optimized using variational inference.

We introduce our framework, GenDLN, where we retain the stacked LLM structure and joint prompt optimization introduced in DLN, but replace the variational inference-based optimization with a Genetic Algorithm (GA) (Fig. 2). The ad-

064

042

043



Figure 2: High-level GenDLN Optimization Framework. Initialization starts from a bank of manual prompts, with optional LLM augmentation. Selection, crossover, and mutation follow the chosen strategies. P: starting population. P': population post-genetic operators.

vantage of using a GA is the ability to explore a large search space and end up with a large pool of candidate prompts. We apply our framework to legal document classification, which aims to categorize legal documents into predefined classes, specifically, Terms of Service (ToS) classification. Also known as Terms and Conditions or Terms of Use, ToS are legal agreements between a service provider and its users, sometimes employing deliberately confusing language (Yerby and Vaughn, 2022), or featuring unfair clauses to users (Loos and Luzak, 2021). Due to ToS length and complexity, users often accept them without fully reading them. To that end, automated unfair clause detection allows consumers to better assess ToS in less than the 45 minutes required to completely read an average ToS agreement (Obar and Oeldorf-Hirsch, 2020). We run further tests on the Microsoft Research Paraphrase Corpus (MRPC) for paraphrase detection.

067

071

073

077

094

Our contributions include a GA framework that successfully improves a population of prompt pairs for classification across several runs and parameter sets, performing in line with state-of-the-art prompt optimization methods. More importantly, our main contribution is an efficient, parameter-rich, LLMbased genetic algorithm framework for text editing that tackles several problems of applying GAs to prompt optimization, including the bottleneck of using API calls for prompt scoring and the additional overheads and limitations imposed by commercial LLM providers. GenDLN can be used by teams with limited resources to quickly generate a pool of optimized prompts for a given task.

### 2 Background

#### 2.1 Prompt Optimization

Prompt optimization is the process of systematically refining or designing the textual instructions (prompts) that guide a Large Language Model toward producing higher-quality, task-specific outputs. Various prompt optimization methods have emerged in recent years. Reflection-based frameworks (Pryzant et al., 2023; Ma et al., 2024) collect error feedback or "textual gradients" from LLM output, then edit prompts accordingly, while editbased approaches (Prasad et al., 2023) iteratively rewrite instructions using operations such as paraphrasing and swapping. Some methods take a meta-prompts approach (Yang et al., 2024), dynamically updating instructions based on historical performance. Additionally, evolutionary algorithm-driven solutions (Guo et al., 2024) simulate natural selection and evolve a population of prompts across generations. All these methods share the same objective: balancing exploration of different prompt variations with exploiting the most promising edits in order to improve the LLM's ability to follow instructions across a range of tasks. In the next sections, we introduce the prompt optimization background used in GenDLN.

101

102

103

104

105

106

107

108

109

110

111

112

113

114

115

116

117

118

119

120

121

122

123

124

125

126

127

128

129

130

131

132

133

134

135

136

137

138

139

140

141

142

143

144

145

146

147

148

#### 2.2 The Stacked LLM

Chaining, stacking, and joining different LLMs has been increasingly explored (Lu et al., 2024; Villarreal-Haro et al., 2024; Ferreira et al., 2024) and shown to perform well across domains for various use cases. The stacked LLM, where outputs from one LLM serve as inputs for another, has proven useful for decomposing complex tasks. One LLM processes raw input, generating intermediate representations or insights; another interprets these representations to complete tasks (classification, reasoning, decision-making). This decomposition boosts accuracy and interpretability (Zhang et al., 2021b) and enhances performance through specialization. Since LLMs excel at when narrowly prompted, this division of labor reduces individual LLM loads and improves result quality (Krishnamurthy et al., 2023). It also allows greater flexibility and modularity in solution design (Khot et al., 2023) while enhancing interpretability, as intermediate outputs clarify reasoning steps (Proca et al., 2024), crucial in fields where black-box decisionmaking is unsuitable, such as law. Lastly, this stacked paradigm mirrors human inference ("First,

244

245

246

247

248

249

201

202

analyze and interpret. Second, draw conclusions and decide" (Correa et al., 2023)). Regardless of the optimization method, stacked LLM architectures offer a clear advantage.

149

150

151

152

153

154

155

156

158

160

161

162

163

164

165

166

167

168

170

171

172

173

174

175

176

177

179

181

183

184

186

188

189

191

192

193

194

195

196

197

199

Sordoni et al. (2023) introduced DLNs as a prompt optimization technique leveraging chained LLM calls. Like other prompt optimization methods, the goal is to use frozen-weight LLMs for inference while refining input prompts for better results. They present two models: DLN-1 (singlelayer) and DLN-2 (two-layer), treating LLMs as stochastic language layers with learnable natural language prompts as parameters. In DLN-2, the first layer's output is considered a latent variable requiring inference, while prompts are learned as parameters of the generative distribution. It employs variational inference for joint prompt optimization in the stacked LLM structure. The advantage of the stacked LLM in DLN is the ability to perform multi-step reasoning through the chaining of prompts and outputs. Similar to the stacked DLN-2 framework, our approach jointly optimizes a prompt pair  $(p_1, p_2)$  for classification, where the scoring function depends on classification metrics. We use the term "DLN" to refer to a two-layer deep neural network (DLN-2). Fig. 1 illustrates GenDLN's prompt pair evaluation.

Unlike DLN, which employs variational inference, we use an LLM-assisted GA for optimization. While LLMs exhibit reasoning-like behavior, research on their stability is mixed, showing high randomness and incoherence (Ma et al., 2024). To mitigate this, we rely on a heuristic optimization strategy (GA) with a reasoning-based evaluation step (DLN) to balance these issues while exploring a larger search space.

#### 2.3 Genetic Algorithms

Genetic Algorithms (GAs) are a class of Evolutionary Algorithms (EAs), global stochastic optimization techniques inspired by Darwin's Theory of Evolution and Natural Selection. They iteratively evolve a "population" of candidate solutions toward the fittest, where the best individual represents the optimal solution (Holland, 1992). Evolutionary approaches excel where traditional methods like gradient descent fail-when the search space is vast, complex, or non-differentiable (Yu and Liu, 2024). Starting with an initial population, candidates are evaluated using a fitness function, with high-fitness individuals more likely to be selected for crossover. Crossover combines features from parents to generate offspring, which serve as new solutions. To maintain diversity, mutations-random occasional changes-are introduced. Repeating this cycle over multiple generations steadily refines solutions, making EAs effective for black-box optimization with minimal system knowledge.

Using GAs for prompt optimization is not new; GAs are proven metaheuristic prompt optimization methods (Pan et al., 2024), with few-shot genetic prompt search surpassing manual tuning (Xu et al., 2022) and evolutionary principles successfully applied to tasks like game comment toxicity classification (Taveekitworachai et al., 2024), preprompt optimization for mathematical reasoning (Videau et al., 2024), Japanese prompting (Tanaka et al., 2023), and emotional analysis (Resendiz and Klinger, 2024). EvoPrompt (Guo et al., 2024) employs LLMs for evolutionary operations like crossover and mutation while EAs guide optimization. The framework implements only one type of selection, crossover, and mutation, all executed by LLMs based on generic instructions, using both manual and LLM-generated initial populations. Our approach, GenDLN (Fig. 2), performs joint prompt-pair optimization instead of single prompt optimization, introduces multiple selection, crossover, and mutation strategies, and implements a richer parameter pool for the GA.

#### 3 Methodology

GenDLN is a multi-objective, steady-state, hybrid genetic algorithm. More details on GenDLN's GA characterization can be found in Appendix A. In this section, we outline the 5 steps of the GA prompt optimization lifecyle in GenDLN (Fig. 2).

**Initialization** (3.1): An initial population of prompt pairs  $(p_1, p_2)$  is sampled from a predefined prompt bank, with optional augmentation.

**Fitness Computation** (3.2): Each individual is scored based on classification metrics by running it through the DLN.

**Selection** (3.3): Individuals are chosen based on fitness using various implemented strategies. **Genetic Operators** (3.4):

**Crossover** (3.4.1): Combines two parents to generate semantically valid offspring.

**Mutation** (3.4.2): Introduces controlled variations to explore new solutions.

**Replacement** (3.5): The next generation is formed by selecting the top individuals, and early stop cri-



Figure 3: Efficiency strategies implemented as part of GenDLN. *Not shown:* workspace level rate-limiter that keeps the frequency of API calls below the platform-defined limit.

teria are defined.

250

251

257

258

262

265

269

272

273

277

#### 3.1 Population Initialization

A population P is a set of individuals. Chromosome encoding refers to how an individual is represented. Each individual I is a prompt pair  $(p_1, p_2)$ , where  $p_1$  is the first-layer prompt for added context and  $p_2$  is the second-layer prompt for classification.

For a population of size N, the initial population consists of N pairs  $(p_1, p_2)$  sampled from a predefined prompt bank, where example prompts are manually added. If the selected size exceeds the available prompts, a Population Initialization LLM optionally generates additional diverse prompts using the prompt bank as examples. Details are in Appendix B.

#### 3.2 Fitness Function / Scoring

The fitness of a prompt pair is computed as a weighted sum of classification metrics, including accuracy and F1 scores, using a multi-objective scoring approach. Fitness is evaluated by running the individual through the DLN (Fig. 1) and comparing predicted labels  $\hat{y}$  to ground truth y. Metric weights are configurable per GA run to reflect different classification goals. Invalid individuals (e.g., with empty prompts) are assigned a fitness of -1 to avoid propagation. Additional fitness implementation and system prompt details for output specification are in Appendix C and E.

278 Rate-Limiting Step: DLN Evaluation The bot279 tleneck in GenDLN is the evaluation of individuals
280 through the DLN, which requires two sequential
281 API calls per data point. Since genetic algorithms
282 require exploring large populations over many gen-

erations, and given the need to use larger models due to the limitations of using smaller ones for prompt optimization (Zhang et al., 2024b), this becomes both time- and cost-prohibitive. To address this, we implement fitness caching, rate limiting, and concurrency across two levels (Fig. 3). First, at the population level (above the dotted line), individuals are evaluated in parallel across w independent workspaces (each representing a compute node with its API key), creating w jobs J that evaluate subsets of the population. Second, within each job J (below the dotted line), the dataset is split into batches. Each batch is then classified using two API calls (one per DLN layer) instead of two calls per data point. Before evaluating a prompt pair  $(p_1, p_2)$ , we check for its presence in a persistent fitness cache. If found, stored metrics are reused, avoiding an expensive DLN evaluation. These optimizations, for a dataset of size 100, increase throughput from  $\approx 18$  to  $\approx 300$  individuals/hour on an 8-core machine, a 16-fold improvement. More details on efficiency strategies and throughput computation are in Appendix F.

284

285

289

291

292

293

294

295

296

297

300

301

302

303

304

305

306

307

308

309

310

311

312

313

314

315

316

317

318

319

320

321

322

323

324

325

326

327

328

329

330

331

332

333

#### 3.3 Selection

Selection can be considered the driving force of the GA; it determines which individuals from the current population will potentially undergo mutation and crossover (and conversely, which members of the current population are discarded), usually based on some function of the individual's fitness. The key is guiding the evolutionary process towards better solutions by preferentially selecting for higher fitness while maintaining population diversity, which is essential to avoid premature convergence. Selection pressure refers to the degree to which individuals with higher fitness are favored during the selection process and directly influences the balance of exploration and exploitation. Higher selection pressure increases the likelihood that fitter individuals will be chosen to pass on their genes, favoring exploitation. This may result in more rapid convergence but also premature convergence if diversity is lost too quickly. Conversely, lower selection pressure allows for a more diverse set of individuals to be selected, favoring exploration but potentially slowing down convergence (Haasdijk and Heinerman, 2018).

The choice of selection strategy is a parameter in GenDLN. We implemented most of the commonly used GA selection strategies, where each has distinct characteristics and influences the

algorithm's selection pressure and, thus, explo-334 ration/exploitation. Selection is the only genetic 335 operator in GenDLN that is not fully or partially 336 LLM-assisted. We implement Random Selection (lack of selection strategy used for comparison purposes), Roulette Wheel Selection (Holland, 1975), 339 Tournament Selection (Miller et al., 1995), Rank-Based Selection (Baker, 2014), Stochastic Uni-341 versal Sampling (SUS) (Baker et al., 1987), and Steady-State Selection. More details on each strat-343 egy's exploration-exploitation balance and implementation can be found in Appendix G. 345

**Preprocessing and Elitism** Before applying any selection method, an optional parameter elitism (k) is used to directly preserve the top k individuals with the highest fitness scores. This ensures that the best-performing solutions are not lost due to stochastic selection effects. For fitness score ties, indices are shuffled, and ties are broken randomly. When  $k \neq 0$ , the individuals are ranked by fitness, and the top k elites are selected for direct inclusion in the next generation. The remaining individuals undergo selection according to the chosen strategy.

347

348

353

354

357

362

365

367

368

371

372

377

379

383

#### **3.4** Genetic Operators in the Textual Space

Since our chromosome is encoded as a tuple of two strings, applying typical crossover/mutation strategies presents challenges. Crossover and mutation are usually performed on bitstrings, numeric vectors, or structured representations of individuals, often following deterministic rules involving slicing, recombining, or editing genes based on strict positional encoding, which is straightforward for bitstring and numeric chromosomes. In the textual space, this is more complex. We discuss these considerations in Appendix H. Work on grammaticallybased genetic programming (Whigham et al., 1995) for creating computer programs has shown the complexity of this task, even in code and query optimization (arguably easier to tokenize than natural language but still sufficiently character- and tokensensitive) (Whigham, 1995).

Research on genetic programming for natural language generation emphasizes the importance of maintaining semantic and syntactic coherence (Araujo, 2020). Thus, we leverage LLMs' ability to dynamically interpret, generate, and refine text as crossover and mutation operators, with prompts passed to an LLM. The response is parsed using regex-based JSON extraction to obtain children in crossover and the mutated prompt in mutation, with a fallback for invalid responses, detailed in Appendix D. Although we have iteratively tested various Mutation and Crossover prompts across different LLMs and included stable ones in GenDLN, these operations remain dependent on LLM responses, with results varying by model and temperature.

384

385

386

388

389

390

391

392

393

394

395

396

397

398

399

400

401

402

403

404

405

406

407

408

409

410

411

412

413

414

415

416

417

418

419

420

421

422

423

424

425

426

427

428

429

430

431

#### 3.4.1 Crossover

We define a set of crossover strategies to allow different levels of exploration and exploitation. The LLM is crucial in ensuring that the offspring are grammatically valid, structurally coherent, and meaningful. We implement 5 strategies: **Single-Point**, **Two-Point**, **Semantic Blending**, **Phrase Swapping**, and **Token-Level** crossover. Details about their implementation and behavior can be found in Appendix I. Crossover is applied to individuals with a user-defined "crossover rate"  $C_r$ , the probability of an individual getting picked to participate in a crossover, and each crossover operation between 2 parents yields 2 children.

#### 3.4.2 Mutation

Much like crossover, we define a set of different mutation strategies leveraging LLMs. The challenge with mutation is the necessity of "limiting" the edits to only a portion of the prompt, as mutation is typically used to introduce comparatively small changes to the chromosome with a userdefined mutation rate  $M_r$ . The goal of mutation is to introduce controlled diversity into the population while maintaining the semantic and syntactic coherence of the prompts. We implement 8 different mutation strategies (Random, Swap, Scramble, Inversion, Deletion, Insertion, Semantic, and Syntactic) with different editing modalities, whose details and prompts can be found in Appendix J.  $M_r$  sets the probability of a "gene" (in our case, a prompt is a gene) to undergo mutation. A "mutate elites" boolean parameter can be used to protect elites from mutation when elitism  $k \neq 0$ . Our choice of strategies and corresponding prompts for both crossover and mutation were made based on our experience and trial and error during the framework's development. It allows for easy editing/extension to include more crossover/mutation types and different prompts. Invalid responses are dealt with using the same retry-fallback mechanism.

#### 3.5 Replacement and Termination

After mutation and crossover, the fitness of the resulting population (now containing approximately  $(N+C_r*N)$  individuals) is calculated, and the top N individuals are the final population of the current generation, with the fittest one being declared the "best in generation."

A GA run is defined for a specific number of generations, but optional stopping criteria can be set, and the GA run will terminate when one of them is met. A "fitness goal" can end the run when the best individual achieves a fitness score equal to or greater than the goal, and a maximum number of stagnant generations S can be set to prematurely terminate the run if the best individual's fitness doesn't improve for S consecutive generations. Otherwise, the GA runs for the predetermined number of generations.

3.6 Logging and Post-Processing

GenDLN features a modular, detailed log structure that allows full retracing of any run. It logs abstractions like best/worst individuals per generation, average metrics, and genetic operator details, alongside full and extracted LLM responses. System details and runtime are also recorded. The output and logging structure is detailed in Appendix K. While implemented in Python, we provide R scripts for post-analysis and extensive GA lifecycle plotting. GenDLN is open source and easily extensible. Additional plots and reproducibility notes are in Appendix L and N.

The following sections describe experiments for binary and multi-label ToS classification on the CLAUDETTE dataset, and binary paraphrase detection on MRPC.

#### 4 Datasets

### 4.1 CLAUDETTE

The CLAUDETTE dataset (Lippi et al., 2019) focuses on Terms of Service agreements from major online platforms, identifying potentially unfair clauses. It includes 50 contracts from providers like Dropbox, Spotify, Facebook, and Amazon, totaling 12,011 sentences, with 1,032 labeled as potentially unfair. Each document is annotated for two classification tasks: binary classification (fair vs. unfair) and multi-label classification, where unfair sentences receive one or more unfairness categories. These include Arbitration, Unilateral change, Content Removal, Jurisdiction, Choice of Law, Limitation of Liability, Unilateral termination, and Contract binding upon usage. Experts manually labeled sentences based on EU consumer law guidelines and court rulings. The dataset is imbalanced across both tasks. For our experiments, we split the data into train, test, and validation sets. LegalBERT and SVM baselines use the full training set, while prompt optimization baselines (OPRO and GrIPS) and our method use a balanced subset of 100 samples per task. A 1000-sample test set is used for evaluation. 

#### 4.2 Microsoft Research Paraphrase Corpus

The Microsoft Research Paraphrase Corpus (MRPC) (Dolan and Brockett, 2005) is a standard benchmark for sentence-level semantic equivalence. It contains 5,801 sentence pairs from news sources, labeled for binary paraphrase detection. We chose MRPC to evaluate GenDLN on a more general, smaller dataset that may not suit fine-tuning or traditional, non-prompt optimization methods. Despite its popularity, MRPC includes formatting artifacts that complicate its use in output-constrained LLM pipelines. We therefore created an *LLM-safe* version via two key preprocessing steps:

**Quote sterilization:** All quote characters (e.g., smart, curly, raw double quotes) were replaced with a Unicode-safe symbol to prevent JSON serialization errors. Mismatched or dangling quotes were manually corrected.

**Trigger filtering:** We removed examples containing high-risk commercial LLM trigger terms.

Our LLM-safe version (See Appendix P for details) preserves task structure and label distribution while ensuring compatibility with LLM-based classification. For experiments, we used 100 balanced training samples and 1000 stratified test samples.

#### Baselines

We compare our approach to both state-of-the-art and classical prompt optimization methods. Optimization by PROmpting (OPRO) (Yang et al., 2024) iteratively refines prompt instructions using an LLM. It uses a meta-prompt containing a problem description, top-performing instructions, and task examples to guide the LLM in generating and evaluating new prompts. Since Legal-BERT performs well in legal NLP classification (Chalkidis et al., 2020), we fine-tuned it on the full training set for all tasks. Our SVM baseline uses TF-IDF vectorization and is trained separately for each task



Figure 4: Metrics (best individual) and average fitness (population) for best CLAUDETTE multi-label run in Table 1.

on the full training set. The other baselines use the same data splits as our approach. The most comparable method to ours is GrIPS (Gradientfree, Edit-based Instruction Search) (Prasad et al., 2023), which edits prompts via deletion, addition, and word swapping, as well as paraphrasing using another LLM. Unlike our approach, it uses simple edit operations and selects top prompts deterministically, without stochastic operators like mutation or crossover.

#### 6 Results and Discussion

We ran over 110 GenDLN executions on CLAUDETTE with various parameter sets across both tasks (binary and multi-label), and around 35 on MRPC. All runs draw from the same set of 10 binary and 10 multi-label manual prompts for CLAUDETTE, and 25 for MRPC, shown in Appendix B, Tables 3–7.

Table 1 lists the runs yielding the bestperforming prompts across the different parameter sets we tried, selected based on Macro F1 performance on the test set. The full prompts for the runs are in Appendix M, Tables 11, 12 and 13. Common parameters for all reported runs: k = 1; no elite mutation; fitness function: 0.2 \* (accuracy) +0.4 \* (macro avg. F1) + 0.4 \* (weighted avg. F1).We use Mistral Large ("mistral-large-2411", 123B parameters), with LLM temperatures for initialization, crossover, and mutation set to 0.7.

Fig. 4 shows the best non-stagnating multilabel CLAUDETTE run (Table 1). Interestingly, it used an insertion mutation strategy, leading to longer prompts, suggesting insertion is exploratorysupported by the diversity plot 7 in Appendix N, which shows a consistently diverse population after the first few generations. While shorter prompts of-



Figure 5: Metrics (best individual) and average fitness (population) for best CLAUDETTE binary run in Table 1. Individual metric lines overlap in the binary case.



Figure 6: Metrics (best individual) and average fitness (population) for best MRPC binary run in Table 1.

ten yield better results (Brown et al., 2020), this run did not early-stop, and could improve with more generations.

566

567

568

570

571

572

573

574

575

576

577

578

579

580

581

582

583

584

585

586

588

Fig. 5 presents metrics for the best binary CLAUDETTE run. Like the multi-label case, we observe stable convergence and fitness improvements across generations. Table 1 lists the best binary run parameters. Unlike multi-label runs, where high-performing prompts were longer, binary runs maintained a more stable prompt length, suggesting structural modifications were more effective than exploratory insertions.

Fig. 6 shows the best MRPC run. MRPC runs resulted in an improvement in accuracy of 6 percentage points on average, with the range of improvement between 3–8 percentage points. Overall, GenDLN consistently improves initial prompts across reasonable parameter settings and remains stable over diverse configurations, and this consistency holds across both datasets. Appendix M includes additional selected runs, parameters, best prompts, and results. Appendix N contains further plots on metrics, convergence, diversity, and

561

565

530

531

	Task	Fitness	Per	formance (T	Test)		GA	A Para	meters				Early
			Acc.	Macro F1	W. F1	Sel.	Cross.	$C_r$	Mut.	$M_r$	Pop.	Gen.	Stop
UDETTE	Binary	0.879	0.79	0.652	0.826	Rank	Sem. Blend	0.8	Semantic	0.2	10	16	Yes
CLA	Multi	0.938	0.825	0.862	0.856	Rank	Phrase Swap	0.85	Insertion	0.3	30	30	No
MRPC	Binary	0.849	0.813	0.796	0.816	Steady-State	Single Point	0.85	Semantic	0.20	30	16	Yes

Table 1: Best GenDLN runs across tasks and datasets. Dataset label is shown in first column. GA Parameters include Selection, Crossover and Mutation types, Population and Generation size, crossover rate  $C_r$  and mutation rate  $M_r$ . Early stop indicates that the run stopped early due to stagnation. W. F1: Weighted F1 score.

		CLAUDETTE					MRPC			
	Binary			Multi			Binary			
	Acc.	Macro F1	W. F1	Acc.	Macro F1	W. F1	Acc.	Macro F1	W. F1	
GenDLN	0.79	0.65	0.83	0.83	0.86	0.86	0.81	0.80	0.82	
OPRO	0.80	0.64	0.83	0.71	0.84	0.84	0.80	0.77	0.80	
(Legal-)BERT*	0.94	0.85	0.94	0.97	0.91	0.91	0.80	0.78	0.80	
SVM TF-IDF	0.93	0.79	0.93	0.77	0.86	0.86	0.70	0.59	0.66	
GRIPS	0.82	0.45	0.85	0.94	0.82	0.82	0.79	0.76	0.79	

Table 2: Test set performance comparison of baseline optimizers across datasets. W. F1: Weighted F1 score. \*BERT was used for MRPC, Legal-BERT was used for CLAUDETTE.

similarity for our best runs (Tables 8, 9, and 10 in Appendix M). We also conduct an ablation study on the best runs, re-running them with "random selection" to isolate selection impact. As expected, ablation results show flatlined metrics, confirming that removing selection pressure collapses the GA into random search (details in Appendix O).

590

592

593

594

611

Our results align with expected GA behavior. All our runs had a maximum population size and generation size of 30, which is the bare-minimum, exploratory number for GA convergence. Rather 599 600 than declaring "optimal" parameter sets for specific tasks, we demonstrate that GenDLN converges across diverse settings, tasks, and datasets. Additionally, Table 2 highlights GenDLN's strong performance against state-of-the-art baselines. In CLAUDETTE binary classification, GenDLN outperforms OPRO and GRIPS in macro F1-score (our 606 prioritized metric due to dataset imbalance). Although LegalBERT and SVM reach the highest scores overall, they rely on full dataset fine-tuning 610 and are not viable for prompt-based few-shot settings. In contrast, GenDLN consistently improves across reasonable parameter configurations using 612 613 only 100 examples. This is especially notable for MRPC, which, unlike CLAUDETTE, does not re-614 quire domain specificity. GenDLN achieves the 615 overall best performance, and is in line with the 616 highest few-shot F1 benchmark of 78.3 in the liter-617

ature reported by Zhang et al. (2021a). For multilabel ToS classification, GenDLN also delivers strong macro and weighted F1 scores, outperforming OPRO and GRIPS in both and surpassing SVM in accuracy, demonstrating its ability to optimize prompt pairs effectively without requiring extensive model adaptation.

618

619

620

621

622

623

624

625

626

627

628

629

630

631

632

633

634

635

636

637

638

639

640

641

642

643

644

645

#### 7 Conclusion

We introduce GenDLN, an efficient evolutionary algorithm-based framework for joint prompt optimization using a stacked LLM architecture. Our approach successfully refines populations of prompt pairs, achieving strong performance on ToS classification and paraphrase detection, in line with baselines such as OPRO and GrIPS on the CLAUDETTE dataset while remaining relatively cost and computationally efficient compared to traditional GA implementations. Through its efficiency strategies at several levels, GenDLN leverages commercial API free tiers to optimize prompt pairs at no cost. This enables resourcelimited teams to use commercial LLMs for EAbased prompt optimization, as applied to welldefined tasks. Our findings highlight the potential of evolutionary strategies as a scalable alternative to traditional prompt engineering and finetuning, paving the way for more accessible and cost-effective LLM-driven classification methods.

743

744

745

746

697

#### 8 Limitations

647

664

669

671

672

676

679

686

Given its reliance on classification based on extraction from an LLM response, the fitness function is subject to model biases and can be influenced by factors such as dataset quality, prompt structure, and stochastic behavior of LLMs. Consequently, fitness scores in this framework serve as an approximation of the true generalization ability of candidate solutions.

Another limitation is the necessary use of the cache, which may lead to bias, but is essential if using commercial LLMs.

Moreover, our framework is limited to tasks/problems where it is possible to encode a solution as a semi-structured, multi-dimensional individual that lends itself to crossover and mutation, and can be assessed by a fitness function. For reasoning/analysis tasks, especially those of a legal nature, the suitability of a solution may be less straightforward to encode and evaluate. Such tasks would require looking at a solution as a multistep task (possibly using more DLN layers and a learned-heuristic approach), such as the work done by Chen et al. (2024).

Additionally, due to the modular logging structure, it is possible to run genetic operators individually and post-process their data. As such, it would be interesting to look at the use of LLMs as genetic operators more closely and examine how they compare to the established stochastic methods, and the bias and differences among different LLMs, temperatures, and parameters.

LLMs are known to sometimes suffer from uncontrolled bias (Bender et al., 2021; Li, 2023; Gallegos et al., 2024). In the context of GenDLN, this may lead to search space restriction due to trigger word sensitivity (Zhao et al., 2024), pretraining bias (Mina et al., 2025), and over-optimization bias (since LLMs are trained to minimize loss on text generation rather than maximize diversity). We have observed anecdotal evidence and instances of the above occuring, but this needs formal exploration.

Furthermore, we do not vary the LLMs and temperature parameters across our different runs. Ideally, instead of relying on the same LLM for all GA operations, different models for mutation, crossover, and evaluation can be used. This approach would introduce flexibility and attempt to reduce systemic bias. Since mutation requires diversity, and a model that introduces novelty, an open model would allow unfiltered, exploratory mutations. Crossover, on the other hand, requires consistency and meaning preservation, and an instructiontuned LLM would be more suitable. For the DLN, a task-specific fine-tuned model would be more reliable for consistent classification.

Moreover, it is important to mention that unlike methods that optimize prompts based on error feedback, GenDLN does not "learn" the dataset in the traditional sense. Due to its reliance on competition and exploration-driven evolution, it shows adaptive improvement, and optimizes prompt pairs for classification with the specific target LLM model used for optimization. This is in line with expected EA behavior. For this reason, specific signals from the dataset will not necessarily make their way to the optimized prompts, and any learning is implicit and general, rather than dataset-specific.

Importantly, we include strong system prompts (based on trial and error) to supplement our optimized prompt pairs. Recent work has explored optimizing system prompts (Zhang et al., 2024a); a development of the idea would be to refine our chromosome encoding to include system prompts. This would make the chromosome carry more than a couple of genes, which is typically the case in GAs.

In addition, due to time constraints, we were not able to run all possible/plausible parameter set combinations. For that same reason, we did not re-run our GA with the same parameter sets, with different random seeds, as is usually done and rely on the high stability (consistent improvement across different parameters sets, tasks, and datasets) of our framework.

Lastly, we only tried one specific fitness score combination, and all our experiments had elitism = 1. We welcome any effort to extend the framework, explore more parameter combinations, and/or formalize parameter exploration for GenDLN through grid search or other techniques.

#### References

- E. Alba and M. Tomassini. 2002. Parallelism and evolutionary algorithms. *IEEE Transactions on Evolutionary Computation*, 6(5):443–462.
- Lourdes Araujo. 2020. Genetic programming for natural language processing. *Genetic Programming and Evolvable Machines*, 21.
- James E Baker et al. 1987. Reducing bias and inefficiency in the selection algorithm.

- 747 748 751 755 756
- 757 761 763 770 772 774 775 782
- 789
- 791 794 796

- 797

- 801

- James Edward Baker. 2014. Adaptive selection methods for genetic algorithms. In Proceedings of the first international conference on genetic algorithms and their applications, pages 101–106. Psychology Press.
- Emily M. Bender, Timnit Gebru, Angelina McMillan-Major, and Shmargaret Shmitchell. 2021. On the dangers of stochastic parrots: Can language models be too big?.
- Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Subbiah, Jared Kaplan, Prafulla Dhariwal, Arvind Neelakantan, Pranav Shyam, Girish Sastry, Amanda Askell, Sandhini Agarwal, Ariel Herbert-Voss, Gretchen Krueger, Tom Henighan, Rewon Child, Aditya Ramesh, Daniel M. Ziegler, Jeffrey Wu, Clemens Winter, Christopher Hesse, Mark Chen, Eric Sigler, Mateusz Litwin, Scott Gray, Benjamin Chess, Jack Clark, Christopher Berner, Sam McCandlish, Alec Radford, Ilya Sutskever, and Dario Amodei. 2020. Language models are few-shot learners.
  - Ilias Chalkidis, Manos Fergadiotis, Prodromos Malakasiotis, Nikolaos Aletras, and Ion Androutsopoulos. 2020. LEGAL-BERT: the muppets straight out of law school. CoRR, arXiv:2010.02559.
  - Yongchao Chen, Jacob Arkin, Yilun Hao, Yang Zhang, Nicholas Roy, and Chuchu Fan. 2024. Prompt optimization in multi-step tasks (promst): Integrating human feedback and heuristic-based sampling.
  - Carlos G. Correa, Mark K. Ho, Frederick Callaway, Nathaniel D. Daw, and Thomas L. Griffiths. 2023. Humans decompose tasks by trading off utility and computational cost.
  - William B. Dolan and Chris Brockett. 2005. Automatically constructing a corpus of sentential paraphrases. In Proceedings of the Third International Workshop on Paraphrasing (IWP2005).
  - Agoston E. Eiben and James E. Smith. 2015. Introduction to Evolutionary Computing. Springer.
  - Tarek El-Mihoub, Adrian Hopgood, Lars Nolle, and Alan Battersby. 2006. Hybrid genetic algorithms: A review.
  - Silvan Ferreira, Ivanovitch Silva, and Allan Martins. 2024. Organizing a society of language models: Structures and mechanisms for enhanced collective intelligence. Preprint, arXiv:2405.03825.
  - Isabel O. Gallegos, Ryan A. Rossi, Joe Barrow, Md Mehrab Tanjim, Sungchul Kim, Franck Dernoncourt, Tong Yu, Ruiyi Zhang, and Nesreen K. Ahmed. 2024. Bias and fairness in large language models: A survey. Computational Linguistics, 50(3):1097-1179.
  - Qingyan Guo, Rui Wang, Junliang Guo, Bei Li, Kaitao Song, Xu Tan, Guoqing Liu, Jiang Bian, and Yujiu Yang. 2024. Connecting large language models with evolutionary algorithms yields powerful prompt optimizers. Preprint, arXiv:2309.08532.

Evert Haasdijk and Jacqueline Heinerman. 2018. Quantifying selection pressure.

802

803

804

805

806

808

809

810

811

812

813

814

815

816

817

818

819

820

821

822

823

824

825

826

827

828

829

830

831

832

833

834

835

836

837

838

839

840

841

842

843

844

845

846

847

848

849

850

851

852

853

- Ali Hakimi Parizi, Yuyang Liu, Prudhvi Nokku, Sina Gholamian, and David Emerson. 2023. A comparative study of prompting strategies for legal text classification. In Proceedings of the Natural Legal Language Processing Workshop 2023, pages 258–265, Singapore. Association for Computational Linguistics.
- Peter J. B. Hancock. 1994. An empirical comparison of selection methods in evolutionary algorithms.
- John H Holland. 1975. Adaptation in natural and artificial systems. University of Michigan Press google schola, 2:29-41.
- John H Holland. 1992. Adaptation in natural and artificial systems: an introductory analysis with applications to biology, control, and artificial intelligence.
- Ziwei Ji, Nayeon Lee, Rita Frieske, Tiezheng Yu, Dan Su, Yan Xu, Etsuko Ishii, Ye Jin Bang, Andrea Madotto, and Pascale Fung. 2023. Survey of hallucination in natural language generation. ACM Comput. Surv., 55(12).
- Tushar Khot, Harsh Trivedi, Matthew Finlayson, Yao Fu, Kyle Richardson, Peter Clark, and Ashish Sabharwal. 2023. Decomposed prompting: A modular approach for solving complex tasks. Preprint, arXiv:2210.02406.
- Yamuna Krishnamurthy, Chris Watkins, and Thomas Gaertner. 2023. Improving expert specialization in mixture of experts. Preprint, arXiv:2302.14703.
- Jinqi Lai, Wensheng Gan, Jiayang Wu, Zhenlian Qi, and Philip S Yu. 2024. Large language models in law: A survey.
- Zihao Li. 2023. The dark side of chatgpt: Legal and ethical challenges from stochastic parrots and hallucination. arXiv preprint arXiv:2304.14347.
- Marco Lippi, Przemysław Pałka, Giuseppe Contissa, Francesca Lagioia, Hans-Wolfgang Micklitz, Giovanni Sartor, and Paolo Torroni. 2019. Claudette: an automated detector of potentially unfair clauses in online terms of service. Artificial Intelligence and Law, 27:117-139.
- Pengfei Liu, Weizhe Yuan, Jinlan Fu, Zhengbao Jiang, Hiroaki Hayashi, and Graham Neubig. 2023. Pretrain, prompt, and predict: A systematic survey of prompting methods in natural language processing.
- M. Loos and J. Luzak. 2021. Update the unfair contract terms directive for digital services. PE 676.006.
- Jinliang Lu, Ziliang Pang, Min Xiao, Yaochen Zhu, Rui Xia, and Jiajun Zhang. 2024. Merge, ensemble, and cooperate! a survey on collaborative strategies in the era of large language models. Preprint, arXiv:2407.06089.

962

- Alessandro Sordoni, Xingdi Yuan, Marc-Alexandre 11
- Ruotian Ma, Xiaolei Wang, Xin Zhou, Jian Li, Nan Du, Tao Gui, Qi Zhang, and Xuanjing Huang. 2024. Are large language models good prompt optimizers? *Preprint*, arXiv:2402.02101.
- Brad L Miller, David E Goldberg, et al. 1995. Genetic algorithms, tournament selection, and the effects of noise. Complex systems, 9(3):193-212.

855

858

870

878

879

890

891

892

894

899

900

901

902

903

904

905

906

- Mario Mina, Valle Ruiz-Fernández, Júlia Falcão, Luis Vasquez-Reina, and Aitor Gonzalez-Agirre. 2025. Cognitive biases, task complexity, and result intepretability in large language models. In Proceedings of the 31st International Conference on Computational Linguistics, pages 1767–1784, Abu Dhabi, UAE. Association for Computational Linguistics.
- Jonathan A. Obar and Anne Oeldorf-Hirsch. 2020. The biggest lie on the internet: ignoring the privacy policies and terms of service policies of social networking services. Information, Communication & Society, 23(1):128-147.
- OpenAI. 2024. Gpt-4 technical report. Preprint, arXiv:2303.08774.
- Rui Pan, Shuo Xing, Shizhe Diao, Wenhe Sun, Xiang Liu, Kashun Shum, Renjie Pi, Jipeng Zhang, and Tong Zhang. 2024. Plum: Prompt learning using metaheuristic. Preprint, arXiv:2311.08364.
- Archiki Prasad, Peter Hase, Xiang Zhou, and Mohit Bansal. 2023. Grips: Gradient-free, edit-based instruction search for prompting large language models. *Preprint*, arXiv:2203.07281.
- Alexandra M. Proca, Fernando E. Rosas, Andrea I. Luppi, Daniel Bor, Matthew Crosby, and Pedro A. M. Mediano. 2024. Synergistic information supports modality integration and flexible learning in neural networks solving multiple tasks.
- Reid Pryzant, Dan Iter, Jerry Li, Yin Tat Lee, Chenguang Zhu, and Michael Zeng. 2023. Automatic prompt optimization with "gradient descent" and beam search. Preprint, arXiv:2305.03495.
- Yarik Menchaca Resendiz and Roman Klinger. 2024. Mopo: Multi-objective prompt optimization for affective text generation. Preprint, arXiv:2412.12948.
- Laria Reynolds and Kyle McDonell. 2021. Prompt programming for large language models: Beyond the few-shot paradigm. In Extended Abstracts of the 2021 CHI Conference on Human Factors in Computing Systems, CHI EA '21, New York, NY, USA. Association for Computing Machinery.
- Côté, Matheus Pereira, Adam Trischler, Ziang Xiao, Arian Hosseini, Friederike Niedtner, and Nicolas Le Roux. 2023. Joint prompt optimization of stacked llms using variational inference. Preprint, arXiv:2306.12509.

- Hiroto Tanaka, Naoki Mori, and Makoto Okada. 2023. Genetic algorithm for prompt engineering with novel genetic operators.
- Pittawat Taveekitworachai, Febri Abdullah, Mustafa Can Gursesli, Antonio Lanata, Andrea Guazzini, and Ruck Thawonmas. 2024. Prompt evolution through examples for large language models-a case study in game comment toxicity classification. In 2024 IEEE International Workshop on Metrology for Industry 4.0 IoT (MetroInd4.0 *IoT*), pages 22–27.
- Mathurin Videau, Alessandro Leite, Marc Schoenauer, and Olivier Teytaud. 2024. Evolutionary pre-prompt optimization for mathematical reasoning. Preprint, arXiv:2412.04291.
- Kapioma Villarreal-Haro, Fernando Sánchez-Vega, Alejandro Rosales-Pérez, and Adrián Pastor López-Monroy. 2024. Stacked reflective reasoning in large neural language models.
- P.A. Whigham. 1995. A schema theorem for contextfree grammars. In Proceedings of 1995 IEEE International Conference on Evolutionary Computation, volume 1, pages 178-.
- Peter A Whigham et al. 1995. Grammatically-based genetic programming. In Proceedings of the workshop on genetic programming: from theory to real-world applications, volume 16, pages 33-41. Citeseer.
- Aled Williams and Yilun Cai. 2024. Insights into weighted sum sampling approaches for multicriteria decision making problems. Preprint, arXiv:2410.03931.
- Hanwei Xu, Yujun Chen, Yulun Du, Nan Shao, Yanggang Wang, Haiyu Li, and Zhilin Yang. 2022. Gps: Genetic prompt search for efficient few-shot learning. Preprint, arXiv:2210.17041.
- Chengrun Yang, Xuezhi Wang, Yifeng Lu, Hanxiao Liu, Quoc V. Le, Denny Zhou, and Xinyun Chen. 2024. Large language models as optimizers. Preprint, arXiv:2309.03409.
- Johnathan Yerby and Ian Vaughn. 2022. Deliberately confusing language in terms of service and privacy policy agreements.
- He Yu and Jing Liu. 2024. Deep insights into automated optimization with large language models and evolutionary algorithms. Preprint, arXiv:2410.20848.
- Lechen Zhang, Tolga Ergen, Lajanugen Logeswaran, Moontae Lee, and David Jurgens. 2024a. Sprig: Improving large language model performance by system prompt optimization. arXiv preprint arXiv:2410.14826.
- Ningyu Zhang, Luoqiu Li, Xiang Chen, Shumin Deng, Zhen Bi, Chuanqi Tan, Fei Huang, and Huajun Chen. 2021a. Differentiable prompt makes pre-trained language models better few-shot learners. arXiv preprint arXiv:2108.13161.

Tuo Zhang, Jinyue Yuan, and Salman Avestimehr. 2024b. Revisiting opro: The limitations of small-scale llms as optimizers. *arXiv preprint arXiv:2405.10276*.

963

964

965

967

968

969

970

971

974

975

976

977

978

979

982

983

985

991

995

997

1001

1003

1005

1006

1007

1008

1009

1010

1011

1012

1014

- Yi Zhang, Sujay Kumar Jauhar, Julia Kiseleva, Ryen White, and Dan Roth. 2021b. Learning to decompose and organize complex tasks. In Proceedings of the 2021 Conference of the North American Chapter of the Association for Computational Linguistics: Human Language Technologies, pages 2726–2735, Online. Association for Computational Linguistics.
- Shuai Zhao, Meihuizi Jia, Zhongliang Guo, Leilei Gan, XIAOYU XU, Xiaobao Wu, Jie Fu, Feng Yichao, Fengjun Pan, and Anh Tuan Luu. 2024. A survey of recent backdoor attacks and defenses in large language models. *Transactions on Machine Learning Research*.

### A GenDLN: GA Characteristics

GenDLN is a multi-objective, steady-state genetic algorithm (SSGA), whereby only a subset of the population is replaced in each generation, and parents evolve alongside their children (through rolling selection, crossover, and mutation) rather than generating an entirely new population. Also, elitism (keeping the best k solutions unchanged) is implemented as an optional parameter, ensuring that the best individual(s) survive to the next generation. Due to employing LLMs in the population initialization, and the mutation and crossover genetic operators, the framework can also be described as a hybrid genetic algorithm (HGA), where domainspecific methods are integrated into the evolutionary process (El-Mihoub et al., 2006). In our domain, textual prompt optimization, GenDLN uses LLM inference to indirectly optimize the initial population, or yield a "good" mutation or crossover product, as opposed to deterministic bit-wise or function-aided manipulations used in classical GAs. Furthermore, in the fitness evaluation, employing the deep-language network (DLN) to determine the suitability of the solution (prompt pair) also makes use of LLM inference and classification-based fitness to guide the optimization process instead of using a deterministic, mathematical function. Our framework is also a multi-objective GA since we use weighted summing of multiple objectives into a single scalar fitness score (Williams and Cai, 2024).

### **B** Population Initialization

This section provides an overview of the population initialization process for the GA, incorporating structured prompt generation and augmentation techniques. **Overview**The population is initialized using pre-<br/>defined sets of prompts, which serve as the basis1015for generating diverse individuals. These prompts1017are loaded and paired to create an initial pool of<br/>candidates. These "prompt banks" as used in our<br/>experiments are shown in Tables 3 and 4.1020

Handling Population SizeIf the predefined set1021of individuals is smaller than the required population size, additional individuals are generated1022lation size, additional individuals are generated1023through augmentation. This ensures a sufficient1024and varied population.1025

Augmentation Process When augmentation is enabled, additional prompts are created by an LLM based on the existing prompt bank. The process ensures that newly created prompts maintain coherence and contribute to the diversity of the population.

**Prompt Generation Details** The augmentation process is guided by a structured system role and user input specification. The following details outline the LLM prompt construction.

#### System Role

You are an expert prompt generator. Based on a given task description and examples, your goal is to generate a specified number of new prompt pairs.

Each prompt pair consists of two prompts:

Prompt 1: An initial instruction to an LLM, to which the LLM would provide a response. Prompt 2: A follow-up instruction, which will be fed to another LLM along with the prior response. This prompt should relate to the expected response from the first LLM.

The new prompt pairs must be unique but adhere to the task description. Each prompt pair must follow this JSON structure:

{
 "prompt\_1": "The first instruction for an
 LLM, which it will respond to.",
 "prompt\_2": "The second instruction for a
 different LLM, which will be given the
 response to prompt\_1 as input and should
 be related to the expected response."
}

```
User Input
```

1062

1026

1027

1028

1029

1030

1031

1032

1033

1034

1035

1036

1037

1038

1040

1041

1042

1043

1044

1045

1046

1048

1049

1050

1051

1052

1053

1054

1055

1056

1057

1058

1059

1060

Prompt 1	Prompt 2
Summarize the following text.	Based on the description above, is the follow-
	ing fair or unfair?
Analyze the following sentence for any poten-	Based on the analysis above, would a reason-
tial ambiguities or unclear language.	able person find this sentence to be fair or un- fair?
Explain the meaning and implications of the	Considering the explanation above, is this sen-
following sentence in a legal context.	tence likely to be considered fair or unfair by a court of law?
Identify any potential biases or assumptions	Based on the identified biases and assump-
present in the following sentence.	tions, do you think this sentence is fair or un- fair?
Provide a paraphrased version of the following	Considering the paraphrased version above, is
sentence to ensure clarity and understanding.	this sentence fair or unfair in its original form?
Discuss the potential consequences of the fol-	Based on the discussion above, do you think
lowing sentence in a legal or regulatory con- text.	this sentence is fair or unfair in its current formulation?
Evaluate the following sentence for compli-	Considering the evaluation above, is this sen-
ance with relevant laws or regulations.	tence fair or unfair in terms of its compliance with applicable laws?
Interpret the following sentence in the context	Based on the interpretation above, is this sen-
of a specific industry or sector.	tence fair or unfair in its application to this industry or sector?
Highlight any potential areas of concern or	Considering the highlighted areas of concern
controversy surrounding the following sen-	above, do you think this sentence is fair or
tence.	unfair in its current form?
Consider the following sentence in light of	Based on the consideration of case law above,
relevant case law or precedents.	is this sentence fair or unfair in terms of its alignment with established legal principles?

Table 3: CLAUDETTE - Manual binary prompt bank used to initialize every GenDLN binary run.

1063	The task involves assessing the fairness or	some insight, analysis, or interpretation	1082
1064	unfairness of a given sentence using prompt	of the given sentence.	1083
1065	pairs.	The response from the first LLM, along with	1084
1066	The sentence will be taken from the Terms &	prompt_2, will then be passed to a second	1085
1067	Conditions or similar documents of	LLM.	1086
1068	businesses. In this context, fairness or		1087
1069	unfairness is meant in the strictly legal	Prompt_2 acts as the final "classifier" and	1088
1070	sense.	<pre>must induce a fair/unfair classification</pre>	1089
1071	This means determining whether the sentence,	by the second LLM based on the response	1090
1072	in isolation, could be deemed unfair by a	to prompt_1.	1091
1073	lawyer or court of law.		1092
1074	We do not know in advance the nature of the	Here are some examples of prompt pairs for	1093
1075	sentence or the area of law it relates to.	the mentioned task:	1094
1076			1095
1077	For this reason, prompt_1 must be general	user_input += "- Prompt 1: " +	1096
1078	and not tied to any specific scenario,	<pre>individual["prompt_1"] + "\n Prompt 2: "</pre>	1097
1079	law, or jurisdiction.	+ individual["prompt_2"] + "\n"	1098
1080	The sentence will be provided alongside		1099
1081	<pre>prompt_1 to an LLM. Prompt_1 will ask for</pre>	user_input += "\nGenerate " +	1100

```
1101 str(total_needed)
1102 + " additional pairs of prompts."
1103
1104 user_input += "Ensure all new pairs
1105 are distinct from the examples."
```

1107

1108

1109

1110

1111

1112

1113

1114

1115

1116

1117

1118

1119

1120

1121

1122

1123

1124

1125

1126

1127

1128

1129

1130

1131

1132

1133

1134

1135

1136

1137

1138

1139

1140

1141

1142

1143

1144

1145

1146

1147

1148

**Finalization** Once the population reaches the desired size, unique identifiers are assigned to each individual. Logging mechanisms help track the composition of the population, distinguishing between original and augmented individuals.

This implementation supports prompt-based population initialization while maintaining flexibility through structured augmentation and validation mechanisms.

# C Fitness Function

The fitness of a prompt pair is a weighted sum of classification metrics using a multi-objective weighted sum approach.

To compute fitness, the individual is evaluated through the DLN (Fig. 1). The classification results  $\hat{y}$  are compared to real labels y, and raw metrics (accuracy, class precision, recall, F1-score, and aggregate metrics like macro- and weighted-average precision, recall, and F1-score) are output by the DLN. Metric weights in the fitness function are configurable per GA run, allowing adaptation to different classification goals, such as prioritizing class-balanced performance by emphasizing macro and weighted metrics or optimizing for specific classes. The sum of metric weights must equal 1, and the resulting fitness score lies in the [0, 1] range. Invalid individuals (where at least one prompt is empty) are assigned a fitness score of -1 to prevent their propagation, as per the fallback mechanism outlined in the next section.

# D Fallback Mechanism for Invalid LLM responses

In GenDLN, LLMs are employed for mutation, crossover and population initialization. The LLM is instructed to generate responses in a valid JSON format, which is necessary for the extraction of prompts and subsequent processing and evaluation of the individuals. However, there are several reasons why the LLM might fail to produce a valid JSON response, beyond ambiguity in prompt instructions (Liu et al., 2023; Reynolds and Mc-Donell, 2021), which is not the case in GenDLN:

1. Model Limitations and Hallucinations:

LLMs are known to potentially "hallucinate" 1149 or generate outputs that deviate from the ex-1150 pected format, especially when the task in-1151 volves complex constraints or novel combi-1152 nations of concepts (Ji et al., 2023). JSON 1153 generation requires strict adherence to syntax 1154 rules, and any deviation (e.g., missing brack-1155 ets, incorrect key-value pairs) results in an 1156 invalid response. 1157

1158

1159

1160

1161

1162

1163

1164

1165

1166

1167

1168

1169

1170

1171

1172

1173

1174

1175

1176

1177

1178

#### 2. Token Limitations and Truncation:

LLMs have a finite context window, and if the generated response exceeds this limit, it may be truncated. Truncation can lead to incomplete JSON structures, rendering the output invalid. This issue is exacerbated when the response includes nested or lengthy JSON objects (OpenAI, 2024).

#### 3. Stochastic Nature of LLMs:

LLMs are probabilistic models, and their outputs can vary significantly even with identical inputs due to temperature settings and sampling strategies. This stochastic behavior increases the likelihood of generating invalid JSON, especially if the temperature parameter is set too high, encouraging creativity at the expense of consistency (Brown et al., 2020). Although our LLM temperature is 0.7 for all experiments, this does not discount the stochastic effects.

#### 4. Crossing Over Identical Prompts:

Some selection strategies naturally lead to 1179 the presence of the same individual more 1180 than once in the population. Moreover, it is 1181 possible to have individuals with one identi-1182 cal prompt through the natural trajectory of 1183 evolution. Since individuals are paired up 1184 for crossover randomly, the Crossover LLM 1185 might be prompted to crossover two "identi-1186 cal" sentences. In most of these cases, the 1187 LLM outputs an invalid response. This was 1188 a problem for all LLMs we tried, including 1189 Llama-3.1-8B, Llama-70B, Ministral 8B, and 1190 even Mistral Large. We chose not to mitigate 1191 this by explicitly including this case and how 1192 to deal with it in the prompt to the operator 1193 and instead detect this with the fallback mech-1194 anism. 1195

## axis in the convergence plots 31, 32, 33, 35, 36,

## **E** System Prompts

37.

**D.1** 

1196

1197

1198

1199

1200

1201

1204

1205

1208

1209

1210

1211

1212

1213

1214

1215

1216

1217

1218

1219

1221

1222

1223

1225

1226

1227

1228

1229

1230

1231

1232

1233

Fallback Mechanism

To mitigate these issues, we implemented a fall-

back mechanism that retries the operation up to a

specified limit (3 in our experiments). If all retries

fail, an empty string is returned, which is detected

during fitness calculation. The assignment of a fit-

ness score of -1 to such individuals ensures that they are not propagated further in the evolutionary

process, maintaining the integrity of the popula-

tion. This approach aligns with established prac-

tices in evolutionary computation, where invalid or malformed individuals are penalized to prevent their influence on future generations (Eiben and

Smith, 2015) and limit their downstream propaga-

tion. We observe that invalid responses occur quite

frequently, and can be visualizes as "X" on the y-

### E.0.1 System Prompts

GenDLN's DLN implementation includes system prompts in scoring. These specify the input/output format (e.g., JSON), define the task, and may include few-shot examples.

Our approach utilizes four distinct system prompts, corresponding to the two-layer binary and multi-label classification approaches. Each prompt defines the input format, specifies the expected output structure, and ensures consistency in model responses.

All prompts follow a common structure:

- The embedded prompt generated by our GA.
- A description of the input format, including identifiers and sentence text.
- A specification of the expected output format, ensuring valid JSON at the second layer.
- Example inputs and outputs to showcase the expected input and output format.

Few-Shot Examples Each system prompt in-1234 cludes six few-shot examples to guide the model's 1235 responses. For binary classification, we randomly select three fair and three unfair sentences from the 1238 training set, ensuring they are distinct from those used in the optimization task. For multi-label classi-1239 fication, we select six sentences, each representing 1240 a unique class. Additionally, for Layer 2 prompts, 1241 the examples include the feature-enriched output 1242

from Layer 1 to provide a more contextualized in-	1243
put.	1244
This approach ensures a balanced representation	1245
of labels while maintaining consistency across both	1246
classification tasks.	1247
We present the full system prompts in the fol-	1248
lowing sections.	1249
E.1 Binary Classification	1250
E.1.1 System Prompt Layer 1	1251
<pre><prompt 01="" placeholder=""></prompt></pre>	1252
	1253
Input Data	1254
The input data is a dictionary containing	1255
sentences from the CLAUDETTE dataset,	1256
where each entry has:	1257
Key: An identifier	1258
(e.g., "sentence_1", "sentence_2")	1259
Value: The sentence text	1260
	1261
Example Input	1262
{	1263
"sentence_1": "This is the text	1264
representing sentence 1.",	1265
"sentence_2": "This is the text	1266
representing sentence 2."	1267
}	1268
E.1.2 System Prompt Layer 2	1269
<prompt_02_placeholder></prompt_02_placeholder>	1270
	1271
Input Data	1272
The input data is composed of two parts.	1273
The first part ("previous_outputs:")	1274
contains a feature-enriched version	1275
of the user input that has already been	1276
processed by a different LLM and	1277
system prompt. The second part	1278
("sentences_to_classify:") is	1279
a dictionary containing sentences	1280
to classify, where each entry has:	1281
	1282
Key: An identifier	1283
<pre>(e.g., "sentence_1", "sentence_2")</pre>	1284
Value: The sentence text	1285
	1286
Example Input	1287
"previous_outputs": "Feature enriched	1288
version of the	1289
sentences to classify"	1290
"sentences_to_classify":	1291

1292

{

```
"sentence_1": "This is sentence 1.",
1293
                  "sentence_2": "This is sentence 2."
1294
1295
              }
1296
              Output Requirements
              For each sentence, add:
1298
              "classification": "fair" or "unfair".
1299
               "rationale": Explanation highlighting
1300
                             influential words.
1301
1302
              Example Output
1303
1304
              {
                   "sentence_1": {
1305
                       "text": "This is sentence 1.",
1306
                       "classification": "fair",
1307
                    "rationale": "Explain the decision."
1309
                   },
                   "sentence_2": {
1310
                       "text": "This is sentence 2.",
1311
                       "classification": "unfair",
1312
                   "rationale": "Explain the decision."
1313
                   }
              }
1315
1316
1317
              Ensure JSON format is valid!
            E.2 Multi-Label Classification
1318
            E.2.1 System Prompt Layer 1
1319
1320
              <Prompt_01_Placeholder>
1321
              CLAUDETTE Classes:
1322
              - PINC (Pins and Cookies)
1323
              - USE (Usage Restrictions)
1324
1325
              - CR (Content Removal)
              - TER (Termination)
1326
              - LTD (Liability Limitation)
1327
1328
              - A (Arbitration)
              - LAW (Applicable Law)
1329
              - J (Jurisdiction)
1330
              - CH (Changes)
1331
1332
              Input Data:
              A dictionary of "unfair" sentences:
1334
              - Key: Sentence ID (e.g., "sentence_1").
1335
              - Value: The sentence text.
1336
1337
1338
              Example Input:
1339
              {
                  "sentence_1": "We may terminate your
1340
                                account at any time.",
                   "sentence_2": "By using Pinterest,
1342
```

you agree to our	1343
policies."	1344
}	1345
E.2.2 System Prompt Layer 2	1346
<prompt_02_placeholder></prompt_02_placeholder>	1347
	1348
CLAUDETTE Classes:	1349
- PINC, USE, CR, TER, LTD, A, LAW, J, CH	1350
	1351
Input Data:	1352
First Part: "previous_outputs"	1353
- Feature-enriched sentences.	1354
Second Part: "sentences_to_classify"	1355
<ul> <li>Dictionary of sentences.</li> </ul>	1356
	1357
Example Input:	1358
"previous_outputs": "Feature enriched	1359
version"	1360
"sentences_to_classify":	1361
{	1362
"sentence_1": "We may terminate	1363
your Account at any time.",	1364
"sentence_2": "By using Pinterest,	1365
you agree to our policies."	1366
}	1367
Evennla Outnut	1368
	1309
i "sontonco 1": (	1071
"text": "We may terminate	1272
vour account "	1372
"classification": ["TER"]	1373
}	1375
"sentence 2": {	1376
"text": "By using Pinterest	1377
vou agree.".	1378
"classification":	1379
["PINC", "USE"]	1380
}	1381
}	1382
-	1383
Each sentence is classified	1384
into one or more labels.	1385
Ensure JSON validity.	1386
F Efficiency Strategies	1387
F.1 Motivation and Setup	1388
Since we use commercial LLM APIs and GAs re-	1389
quire exploring a vast search space to converge,	1390

running our framework is both cost- and time-

1391

intensive, especially for fitness evaluation. Evalu-1392 ating a prompt pair through the DLN requires two 1393 API calls per data point. For large datasets and 1394 populations (essential for exploration), running the 1395 framework for enough generations becomes too expensive, not to mention the need to test various 1397 parameter sets and the significant trial-and-error 1398 phase inherent to evolutionary optimization. To 1399 mitigate this, we implemented efficiency strategies 1400 at different framework stages. We apply metric 1401 caching, request rate limiters, and concurrency at 1402 two DLN levels (Fig. 3). 1403

#### F.1.1 Metric Caching

1404

1420

1421

1422

1423

1424

1425

1426

1427

1428

1429

1430

1431

1432

1433

1434

1435

1436 1437

1438

1439

1440

1441

As mentioned, running an individual through the 1405 DLN yields a set of classification metrics. In 1406 1407 GenDLN, these raw metrics are cached for every prompt pair to avoid rerunning the evaluation of the 1408 same prompt pair within the same run; we also ex-1409 tend it to avoid rerunning the evaluation of the same 1410 prompt pair for the same LLM-dataset-task combi-1411 nation. The cost savings and speed-up provided by 1412 caching comes at the risk of introducing some bias 1413 (LLM-classification is inherently unstable, and the 1414 same prompt can lead to different responses from 1415 the same LLM). However, this is primarily used 1416 to explore parameter sets, and for suitable, stable 1417 parameter definitions, the GA should eventually be 1418 rerun three times to discount noise. 1419

#### F.1.2 Parallelization

Significant work has been done on parallelizing the execution of GAs (Alba and Tomassini, 2002). For GAs in general, evaluation of an individual is independent, and for GenDLN (DLN classification using prompts  $(p_1, p_2)$ ), this allows population evaluation to be parallelized. To accelerate the prompt optimization process, our framework employs a two-layer parallelization approach, addressing both the evaluation of individual prompt pairs and the internal processing of data batches for each individual.

**Inter-Individual Parallelization** In Fig. 3, the top section (above the dashed line) shows population-level parallelization, our first concurrency layer.

A workspace W is a compute node with an independent API token handling requests. For a w-core machine, w sets of individuals from population P run in parallel across w workspaces, creating w jobs J, each evaluating up to N/w individuals. Rather than processing individuals sequentially, our framework concurrently evaluates sev-1442 eral prompt pairs. This strategy exploits multi-1443 core architectures to significantly reduce the over-1444 all optimization time. By partitioning the popula-1445 tion across multiple execution threads or processes, 1446 each prompt pair can be evaluated independently. 1447 Importantly, each individual maintains its own iso-1448 lated "workspace," meaning that the computational 1449 resources and rate-limiting mechanisms are man-1450 aged on a per-individual basis. 1451

1452

1453

1454

1455

1456

1457

1458

1459

1460

1461

1462

1463

1464

1465

1466

1467

1468

1469

1470

1471

1472

1473

1474

1475

1476

1477

1478

1479

1480

1481

1482

1483

1484

1485

1486

**Intra-Individual Concurrency** The bottom section (Fig. 3) details job J. Within the evaluation of a single prompt pair (job J), further efficiency is gained by concurrently processing the training dataset. We first partition the dataset into multiple batches, then evaluate the prompt pair on these batches concurrently, using 2 API calls (one per DLN layer/prompt) per batch rather than 2 per sentence.

This fine-grained parallelism allows us to aggregate evaluation metrics faster, as each batch is processed in parallel rather than sequentially. The results across individuals and batches are aggregated to determine  $(p_1, p_2)$ 's overall performance, with metrics stored in the cache for future use.

A notable constraint in our setup is the use of an external API that enforces a strict rate limit of one request per second (RPS). To adhere to this limit while still maintaining high throughput, we integrate a rate limiter into our concurrency model. For each prompt pair, the batch-level evaluations are regulated such that API calls are spaced appropriately. Since each individual has its own "workspace," the rate limiting is applied independently per prompt pair. This design ensures that the API is not overwhelmed by simultaneous requests across the entire population while still exploiting concurrency within each evaluation task.

Overall, the combination of inter-individual parallelization and intra-individual concurrency leads to a significant speedup in our prompt optimization process, allowing us to efficiently explore the search space while managing the operational constraints imposed by the external API.

#### F.1.3 Individual Evaluation Throuphput

To quantify the efficiency of our genetic algorithm1487runs, we define the *individual evaluation through-*1488put as the number of individuals evaluated per unit1489of time. Given a genetic algorithm run with G gen-1490erations, a population size of N, a crossover rate of1491

1492

1/10/

149

1496 1497

- 1498
- 1499

1500 1501

1502 1503

- 1504
- 1505
- 1506
- 1507
- 1508
- 1509

# 1510 1511

1512

1513

1514

1515

1516

1517

1518

1519

1520

1521

1522

1523

1525

1526

1527

1528

1529

Random selection is the absence of a selection strat-

**G.1** 

G

as:

across all generations is:

tions by the runtime:

egy. It refers to selecting individuals uniformly at random, irrespective of fitness values. We implement it for use as a baseline for comparison purposes.

 $C_r$ , and a total runtime of T hours, the number of

individuals evaluated per generation is computed

 $N(1 + C_r)$ 

 $G \cdot N(1+C_r)$ 

uals evaluated per hour, we divide the total evalua-

 $\text{Throughput} = \frac{G \cdot N(1 + C_r)}{T}$ 

This metric allows us to compare different ge-

netic algorithm configurations by normalizing their

efficiency in terms of evaluations processed per

hour, thereby accounting for variations in runtime

across different experimental settings.

**Selection Strategies** 

**Random Selection** 

To determine the throughput in terms of individ-

Thus, the total number of individual evaluations

(1)

(2)

(3)

# G.2 Roulette Wheel Selection

Also known as fitness proportionate selection, roulette wheel selection is one of the very first explored GA selection strategies (Holland, 1975). It simulates spinning a wheel where each individual occupies space proportional to its fitness, and selections are made probabilistically (by "spinning" a wheel and selecting the individual the "pointer" lands on). It ensures that individuals with higher fitness have a higher chance of selection, but any individual could potentially be selected. However, if relatively high-fitness individuals dominate early, this may lead to premature convergence. Also, when fitness values are very similar, low selection pressure may lead to stagnation (Hancock, 1994).

1530Tournament SelectionFirst introduced by1531Miller et al. (1995), tournament selection is a simple and widely-used selection strategy. For a tournament size t, it randomly picks t individuals from1533the population, and selects the individual with highest fitness (the "tournament winner") for the next

generation. For a population size N, N tourna-1536 ments are held, with t participants each (if elitism 1537  $k \neq 0, N - k$  tournaments are held). Tournament 1538 selection aims to establish a balance between explo-1539 ration and selection pressure, which can be tuned 1540 with tournament size t. Larger tournaments lead 1541 to stronger selection pressure and lower diversity 1542 (exploitation), while smaller tournament sizes favor 1543 exploration. 1544

1545

1546

1547

1548

1549

1550

1551

1552

1553

1554

1555

1556

1557

1558

1559

1560

1561

1562

1563

1564

1565

1566

1567

1568

1569

1570

1572

1573

**Rank-Based Selection** Conceptually similar to roulette wheel, rank-based selection assigns individuals space on the wheel according to their rank rather than their fitness, where the total space on the wheel is equal to the sum of the ranks. Introduced by Baker (2014), to mitigate scaling issues where individuals in the population have fitness values that are either too extreme (high-fitness outliers would be selected too often in classical roulette), or too similar (if fitness values are too close together, each individual would have roughly the same chance of being selected in classical roulette). Rank selection ensures a linear selection probability distribution which prevents bias towards disproportionately high fitness individuals, while maintaining selection pressure.

Stochastic Universal Sampling (SUS) SUS was introduced by Baker et al. (1987) as an improvement over roulette wheel selection. In this variant, N evenly spaced pointers are assigned to the wheel, on which the individuals occupy space proportional to their fitness values, and N individuals are selected in one go when the wheel is "spun." It ensures a more diverse selection and reduces stochastic noise, but will still suffer from premature convergence in the presence of a high-fitness outlier (if an individual occupies a disproportionately large space on the wheel, several pointers will land on it).

Steady-State Selection Our framework is inher-1574 ently an SSGA due to the way our replacement step 1575 (discussed in a futher section) operates, however, we also implement an explicit steady-state selec-1577 tion strategy for greater flexibility. Steady state 1578 selection requires elitism  $k \neq 0$  or else it will be-1579 have like random selection. In this strategy, the 1580 top k fittest individuals are selected for the next 1581 generation, and N - k are randomly selected from 1582 the remaining individuals to complete the popula-1583 tion. Steady-state selection ensures that only a few individuals are replaced at a time in each genera-1585

diversity.

1589 1590

1617

1618

# H Adapting Chromosomes to the Textual Space - Considerations

tion. Always keeping many elites in the population

may accelerate convergence at the risk of reducing

1591 Although we have encoded the chromosome as a tuple, that does not mean the individual only has 2 1592 genes  $(p_1 \text{ and } p_2)$ . The "suitability" of the solution 1593 depends on unstructured, hard-to-define components or "tokens" within the two text prompts, as 1595 well as hidden "genetic material" in the textual fea-1596 tures of each prompt string. In natural language, 1597 different words, phrases, and clauses hold different 1599 weights in conveying meaning, unlike in structured encoding, where every component's contribution 1600 to the solution's suitability is defined. If classical 1601 strategies were to be applied (slicing the strings 1602 at arbitrary points, editing the characters at arbi-1603 trary indices), this would risk yielding too many 1604 syntactically invalid or semantically nonsensical 1605 prompts. Additionally, words and phrases are in-1607 terdependent (much like real genes), and simple positional swapping and randomized editing may 1608 distort the meaning. In fact, textual meaning can 1609 completely collapse if crossover/mutation is badly 1610 applied, yielding individuals far inferior to their 1611 progenitors, which defeats the purpose. Determin-1612 ing where and how to split/edit text dynamically 1613 while ensuring coherence of results is an inherently 1614 1615 non-deterministic process, contrary to the established concept of crossover and mutation in GAs. 1616

# I Crossover Strategies

We implemented the following strategies:

- 1619Single-PointSelects a single random point in1620each sentence and swaps the latter halves to form1621new sentences.
- 1622**Two-Point**Selects two random points in each1623sentence, swapping alternating segments to form1624new sentences.
- 1625Semantic BlendingBlends the core meaning of1626both parents into two complementary sentences.1627Offspring are not simple recombinations but rather1628semantically fused versions of the inputs.
- 1629Phrase SwappingIdentifies key phrases in each1630parent and swaps them while maintaining grammat-1631ical integrity.

**Token-Level** Swaps individual words or tokens between sentences.

1632

1633

1634

1646

1649

1651

1652

1654

1655

1657

1659

1660

1662

1670

1673

1674

1675

# I.1 Crossover System Prompt

"You are an expert linguist and copywriter, act-<br/>ing similar to how genetic crossover works, but1635in a textual context. Generate two complementary1637sentences as children of the provided parent sen-<br/>tences. Here complementary means that the two1639child sentences must have complementary parts of<br/>the parents, as in genetic crossover. Make sure the<br/>children sentences are wrapped in a JSON-object1642as follows:1643

{"child_1": "child sentence 1",	164
"child_2": "child sentence 2"}	164

The rest of your response can be plain text, but the new sentences must be in a JSON. Both sentences must be grammatically correct and reasonably meaningful."

# I.2 Crossover Strategy Prompts

**Single-Point** "Combine the following two sentences by splitting each at a single random point. The first child should take the first half of the first sentence and the second half of the second sentence. The second child should take the first half of the second sentence and the second half of the first sentence. Ensure both sentences remain coherent and meaningful."

**Two-Point** "Combine the following two sentences by selecting two random points in each sentence. The first child should integrate the segments alternately, starting with the first part of the first sentence. The second child should integrate the remaining segments alternately. Ensure both sentences are coherent and meaningful."

**Semantic Blending** "Blend the following two sentences to create two complementary sentences. Each child should focus on combining the core meaning of both sentences in a unique way. Ensure that both sentences are coherent, meaningful, and distinct from one another."

**Phrase Swapping** "Swap one or more phrases between the following two sentences to create two new sentences. Each child should incorporate phrases from the other parent in a way that creates a coherent and meaningful result."

Token-Level"Swap individual words or tokens1677between the following two sentences to create two1678

- 1681
- 1682 1683
- 1684
- 1685
- 1686 1687
- 1688
- 1689 1690
- 1691
- 1692
- 169
- 1694
- 1695
- 1696
- 1697
- 1698 1699
- 1700
- 1701

1702 1703

- 1704
- 1705
- 1706
- 1707 1708

1709

1711 1712

1713

- 1714 1715
- 1716
- 1717
- 1718
- 1719
- 1720 Phrase Swapping
- 1721Parent 1: "Summarize the following text."1722Parent 2: : "Explain the meaning and1723implications of the following sentence1724in a legal context."

new sentences. Each child should incorporate

words from the other parent in a way that creates a

Below are some selected illustrative crossover ex-

Parent 1: "Summarize the following text."

following sentence in a legal context."

Child 1: "Summarize the following text

Child 2: "Explain the meaning and

Parent 2: "Explain the meaning and

Child 1: "Summarize the meaning and

implications of the following text."

Parent 1: "Summarize the following text."

implications of the following sentence

implications of the following sentence in

Child 2: "Explain the following text in

a concise manner and its potential impact

Parent 1: "Based on the description above,

Parent 2: : "Considering the explanation

above, is this sentence likely to be

considered fair or unfair by a court

Child 1: "Considering the description

above, is the treatment likely to be

considered fair or unfair by a court

Child 2: "Based on the explanation

above, is the sentence likely to be considered fair or unfair in a court

is the following fair or unfair?"

Parent 2: "Explain the meaning

and implications of the

in a legal context."

in a legal context."

a legal context"

**Semantic Blending** 

on the law"

of law?"

of law?"

of law?"

coherent and meaningful result."

I.3 Crossover Examples

amples.

**Single-Point** 

**Two-Point** 

1725Child 1: "Explain the meaning and1726implications of the following summary in

a legal context."1727Child 2: "Summarize the following sentence1728to understand its core message and1729implications."1730

1731

1744

1745

1746

Token-Level

Parent 1: "Based on the description above, 1732 is the following fair or unfair?" 1733 Parent 2: "Considering the explanation 1734 above, is this sentence likely to be 1735 considered fair or unfair by a court 1736 of law?" 1737 Child 1: "Considering the description above, is the following sentence likely 1739 to be considered fair or unfair by a 1740 court of law?" 1741 Child 2: "Based on the explanation above, 1742 is the following sentence likely to be 1743

is the following sentence likely to be considered fair or unfair by a court of law?"

J Mutation Strategies

5	
The following is a summary of the introduced strategies and their intended result.	1747 1748
<b>Random</b> Changes a single word or phrase in the sentence to a synonym or a similar concept.	1749 1750
<b>Swap</b> Swaps existing words or phrases in the sentence to introduce minor structural variation.	1751 1752
<b>Scramble</b> Rearranges the order of words/phrases while maintaining the original meaning.	1753 1754
<b>Inversion</b> Reverses the order of words or phrases in part or all of the sentence.	1755 1756
<b>Deletion</b> Removes a word or phrase from the sentence to create a more concise variation.	1757 1758
<b>Insertion</b> Adds new words or phrases to provide additional context while preserving meaning.	1759 1760
<b>Semantic</b> Rephrases the sentence slightly while keeping the core meaning intact.	1761 1762
<b>Syntactic</b> Alters the sentence structure while preserving the meaning.	1763 1764
J.1 Mutation System Prompt	1765
"You are an expert linguist and copywriter. Make sure the sentence you return is wrapped in a JSON- object as follows:	1766 1767 1768
{"mutated_sentence": "new sentence you generate based on the instruction"}.	1769 1770

## 1775 J.2 Mutation Strategy Prompts

- **Random** "Change only one single word or phrasein the sentence to a synonym or similar concept."
- 1778Swap"Swap two existing words or phrases in the1779sentence."
- 1780Scramble"Rearrange the existing words and/or1781phrases in the sentence with a minimal addition of1782new words."
- 1783Inversion"Invert the order of the existing words1784or phrases in all or part of the sentence."
- 1785Deletion"Delete a word or phrase in the sen-1786tence."
- 1787Insertion"Insert words or phrases in the sen-1788tence that could provide more context/clarity while1789keeping the same base meaning."
- 790 Semantic "Slightly rephrase the sentence."
- 1791Syntactic"Modify the sentence structure of the1792sentence while keeping the same base meaning."

# J.3 Mutation Examples

Below are some selected illustrative mutation examples.

# 1796 Semantic

1793

1794

1795

1805

- 1797 Initial Prompt: "Produce a detailed output
  1798 for each sentence, outlining the reasoning
  1799 for its classification into the most likely
  1800 category."
- 1801Mutated Prompt: "Generate a comprehensive1802output for each sentence, explaining the1803rationale for its categorization into the1804most probable group."

# Insertion

- 1806Initial Prompt: "Interpret each sentence1807and provide a comprehensive rationale for1808its legal classification."
- 1809Mutated Prompt: "Carefully interpret each1810individual sentence within the context of1811the document and provide a comprehensive1812rationale for its specific legal1813classification."

## Random

Initial Prompt: "Summarize the following 1815 text." 1816 Mutated Prompt: "Condense the following 1817 text." 1818

1814

1819

1820

1821

1823

1824

1827

1828

1829

1830

1831

1832

1833

1834

1836

1837

1838

1840

1841

1842

1843

1844

1845

1846

1847

1848

1849

1851

1852

1853

1854

1855

1857

# Swap

Initial Prompt: "Based on the description above, is the following fair or unfair?" Mutated Prompt: : "Based on the description above, is the following unfair or fair?"

## Deletion

Initial Prompt: "Based on the description above, is the following fair or unfair?" Mutated Prompt: "Based on the description, is the following fair or unfair?"

# Scramble

Initial Prompt: "Based on the description above, is the following fair or unfair?" Mutated Prompt: "Is the following fair or unfair, based on the description above?"

# K GenDLN Logging

Every sub-component of GenDLN (fitness calculation, selection, crossover, mutation, replacement, caching) has a dedicated logger and defined structure, and a GA Log (which is the output of the framework), is a structured log of these components. Below we provide the expected output and logger functionality and examples.

The logging system in the Genetic Algorithm (GA) serves as a comprehensive tracking and debugging framework, capturing detailed records of key evolutionary events at multiple levels. It ensures traceability of the entirety of the GA run. The logging structure is hierarchical, with nested loggers handling distinct operations, and a centralized GA logger aggregating all logs.

**Hierarchical Structure of Logging** The logging framework consists of specialized loggers:

- GA Logger The central log for the entire evolutionary process, containing pergeneration records of all key operations.
- **Population Initialization Logger** Tracks how the initial population is created, including augmentation details.
- Selection Logger Records selected individuals, strategy parameters, and elitism effects. 1859

 Crossover Logger – Captures the details of crossover operations, including parentoffspring relationships.

1860

1861

1862

1865

1866

1867

1868

1869

1870

1872

1873

1876

1877

1878

1879

1902

1903

1905

- **Mutation Logger** Stores information on how individuals are mutated, along with mutation types.
- Fitness Logger Logs individual fitness scores and overall generation-level fitness statistics.
- Fitness Cache Logger Tracks cache hits and misses.
- **Replacement Logger** Logs how individuals are retained or replaced in the next generation.
- **Run-Specific Details** Runtime, system specs, configs, and hyperparameters of the GA run are appended to the end of the log.

GA Logger: Centralized Evolution Tracking

Each generation's log entry contains the following:
{
 generations": [

```
1880
                     {"generation_id" : i,
                     "initial_population": [...],
1882
                     "selection_data": [...],
1883
                  "population_after_selection": [...],
1884
                     "crossover_data": [...],
1885
                  "population_after_crossover": [...],
1886
                     "mutation_data": [...],
1887
                  "population_after_mutation": [...],
1888
                     "fitness_data": {...},
1889
                     "replacement_data": [...]
1890
1891
                     },
                      \{\ldots\}, \ldots],
1892
                 "early_stopping":
                     {"status": false, "reason": ""},
1894
                 "runtime": "3.25 minutes",
1895
                 "system_info": {...},
1896
                 "config": {...},
                 "hyperparameters": {...},
1898
                 "ga_log_filename":
1899
                     {"ga_log_date-timestamp.log"}
1900
            }
1901
```

This hierarchical logging system ensures that all operations are transparently recorded, aiding both debugging and performance analysis of the genetic algorithm.

# L Reproducibility

We provide a set of R Scripts that allow the repro-<br/>duction of our results, plots, and analyses. The<br/>scripts are structured to ensure transparency and<br/>ease of replication, and enforce a file path structure<br/>for inputs and outputs.1907<br/>1908

1906

1912

1919

1920

1921

1933

1936

1937

1941

1943

# L.1 Environment Setup

All necessary dependencies are installed and loaded1913at the start of the execution. The required R1914libraries include tidyverse, jsonlite, here,1915purrr, data.table, dplyr, ggplot2, tidyr,1916readr, and stringdist. The script automatically1917installs missing dependencies.1918

# L.2 Data and Directory Structure

The project assumes a structured directory for data storage and result output:

- **Root Directory**: Automatically set to the location of the script. 1922
- Log Directory: Stores raw Genetic Algorithm (GA) log files (output of GenDLN). 1924
- Summary Directory: Contains extracted 1926 metadata and performance summaries. 1927
- Test Directory: Stores test results.
- Output Directory: Stores processed results 1929 and plots. 1930
- **Plot Directory**: Contains visualization outputs. 1931

All necessary directories are created if they do not exist.

# L.3 Processing and Normalization

**Log File Normalization** GA log files are processed into structured formats. Key extracted elements include:

- Initial generation data (fitness scores, raw metrics, attributes).
- Subsequent generation data with performance metrics.
- Total number of completed generations.
- Metadata including runtime, system configuration, hyperparameters, and early stopping conditions.
   1944
   1945
   1946

1948	cessed to extract structured information on:	The script com tracking:
1949 1950	• GA parameters (population size, mutation rate, selection strategy, fitness function).	• Unique in
1951 1952	• Run performance (best fitness scores, accuracy, raw evaluation metrics).	• Similarity tions.
1953 1954	• Execution environment (system specifications, runtime details).	• Levenshte best indiv
1955 1956	L.4 Batch Processing and Summary Generation	<b>Comprehensi</b> bined summary
1957	Aggregating Run Summaries A batch process-	CSV file.
1958	ing script collects metadata from all runs and pro-	
1959	cludes:	M Detailed
1000	ciudos.	M.1 CLAUE
1961	• Number of runs per batch.	The best prom
1962	• Associated test results.	As for mult
1963	• Log files used in the batch.	prompts are in
1964	This process ensures that interrupted runs are	M.2 MRPC
1965	accounted for and test data is linked correctly.	The best prom Table 10 are sh
1966	Appending Notes to Summaries Notes can be	
1967	appended to individual summaries to document	N Detailed
1968	special conditions or anomalies in the runs.	N.1 Metrics
1969	L.5 Analysis and Visualization	The metrics ov
1970	GA Performance Report Each log file is pro-	formance metri
1971	cessed to produce a detailed report that includes:	It is a multi-lin
1972	• Performance metrics across generations (fit-	metric and its
1973	ness scores, accuracy, F1 scores).	represents the g
1974	• Statistical summaries (mean, variance, min,	indicate differe
1975	max values of key metrics).	Higher valu mance. Fluctua
1976	• Evolutionary trends of best and worst individ-	instability or early
1977	uals.	(GA), while a d tion around op
1978	Metric Extraction and Visualization Metrics	ing or stable fit
1979	such as fitness score, accuracy, and F1 scores are ex-	vergence, when
1980	tracted for each generation and visualized to track	score suggests
1981	GA progression.	CLAUDETTE
1982	GA Convergence Analysis The convergence of	are on the left
1983	the GA is visualized by plotting best and worst	binary runs, the
1984	fitness scores across generations.	19, 21.
	2	3

Metadata Extraction Log files are further pro-

1947

Diversity and Similarity of Best Individuals	1985
tracking:	1986
utexnig.	1501
• Unique individuals per generation.	1988
• Similarity of best individuals across genera-	1989
tions.	1990
• Levenshtein and Jaccard similarity scores for best individuals.	1991 1992
Comprehensive Run Summary A final com-	1993
bined summary consolidates all extracted informa-	1994
tion, test results, and log metadata into a structured	1995
CSV file.	1996
M Detailed Desults	1007
WI Detaileu Results	1997
M.1 CLAUDETTE	1998
The best prompts from the top 4 selected binary	1999
runs in Table 8 are shown in Table 11	2000
As for multi-label, results are in Table 9, and	2001
prompts are in Table 12.	2002
M.2 MRPC	2003
The best prompts from the top 4 selected runs in	2004
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13	2004 2005
<ul><li>The best prompts from the top 4 selected runs in Table 10 are shown in Table 13</li><li>N Detailed Plots</li></ul>	2004 2005 2006
<ul> <li>The best prompts from the top 4 selected runs in Table 10 are shown in Table 13</li> <li>N Detailed Plots</li> <li>N.1 Metrics Over Generations</li> </ul>	2004 2005 2006 2007
<ul> <li>The best prompts from the top 4 selected runs in Table 10 are shown in Table 13</li> <li>N Detailed Plots</li> <li>N.1 Metrics Over Generations</li> <li>The metrics over generations plot tracks key per-</li> </ul>	2004 2005 2006 2007 2008
<ul> <li>The best prompts from the top 4 selected runs in Table 10 are shown in Table 13</li> <li>N Detailed Plots</li> <li>N.1 Metrics Over Generations</li> <li>The metrics over generations plot tracks key performance metrics across generations, such as accu-</li> </ul>	2004 2005 2006 2007 2008 2009
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores.	2004 2005 2006 2007 2008 2009 2010
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a	2004 2005 2006 2007 2008 2009 2010 2011
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis	2004 2005 2006 2007 2008 2009 2010 2011 2012
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics.	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics.	2004 2005 2006 2007 2008 2010 2010 2011 2012 2013 2014 2015 2016
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Eluctuations in fitness and accuracy reflect	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabiliza-	2004 2005 2006 2007 2008 2010 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabilization around optimal solutions. A steadily increas-	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabilization around optimal solutions. A steadily increasing or stable fitness score implies progress and con-	2004 2005 2006 2007 2008 2010 2010 2011 2012 2013 2014 2015 2016 2017 2016 2017 2018 2019 2020 2021
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabilization around optimal solutions. A steadily increasing or stable fitness score implies progress and convergence, whereas a volatile or fluctuating fitness	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2022
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key per- formance metrics across generations, such as accu- racy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better perfor- mance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabiliza- tion around optimal solutions. A steadily increas- ing or stable fitness score implies progress and con- vergence, whereas a volatile or fluctuating fitness score suggests ongoing evolution.	2004 2005 2006 2007 2008 2010 2010 2011 2012 2013 2014 2015 2016 2017 2016 2017 2018 2019 2020 2021 2022 2023
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabilization around optimal solutions. A steadily increasing or stable fitness score implies progress and convergence, whereas a volatile or fluctuating fitness score suggests ongoing evolution.	2004 2005 2006 2007 2008 2009 2010 2011 2012 2013 2014 2015 2016 2017 2018 2019 2020 2021 2020 2021 2022 2023
The best prompts from the top 4 selected runs in Table 10 are shown in Table 13 <b>N Detailed Plots</b> <b>N.1 Metrics Over Generations</b> The metrics over generations plot tracks key performance metrics across generations, such as accuracy, fitness score, average fitness, and F1 scores. It is a multi-line plot where each line represents a metric and its trend over generations. The x-axis represents the generation number, while the y-axis represents the value of the metric. Different colors indicate different metrics. Higher values generally indicate better performance. Fluctuations in fitness and accuracy reflect instability or exploration by the genetic algorithm (GA), while a converging trend suggests stabilization around optimal solutions. A steadily increasing or stable fitness score implies progress and convergence, whereas a volatile or fluctuating fitness score suggests ongoing evolution. <b>CLAUDETTE</b> Plots for the top multi-label runs are on the left side of Fig. 7, 9 and 11, 13. For the	2004 2005 2006 2007 2008 2010 2010 2011 2012 2013 2014 2015 2016 2017 2016 2017 2018 2019 2020 2021 2022 2023 2024 2025



Figure 7: CLAUDETTE - Left: plot of metrics and average fitness for best run A in Table 9. Right: Diversity plotting for best multi-label run A in Table 9



Figure 8: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run A in Table 9.

**MRPC** Plots for the top runs are on the left side of Fig. 23, 25 and 27, 29.

#### N.2 Convergence Plot

2028

2032 2033

2034

2035

2036

2037

2039

2040

2041

2046

2047

2050

The convergence plot visualizes how the best and worst individuals change across generations, providing insight into GA optimization progress. This line plot features a dashed blue line representing the best fitness and a dotted red line representing the worst fitness. A shaded region between these lines indicates population fitness spread. The xaxis represents the generation number, and the yaxis represents the fitness score. The best fitness line tracks the top-performing individual in each generation, while the worst fitness line tracks the least-performing individual. A narrowing gap between the two lines indicates that the population is converging toward similar solutions. If the best fitness stagnates early, the algorithm may have prematurely converged to a suboptimal solution. Convergence occurs when the best and worst scores stabilize and remain close together. A wide gap between best and worst scores suggests high diversity in the population. If the worst score is constantly

low, it may indicate poor-quality individuals or unfit solutions. The X on the Y-axis represents a worst individual with an empty prompt, which was detected by the fallback mechanism described in D and assigned a fitness score of -1, not represented in the y-axis scale in order not to skew the graph.

2054

2061

2062

2063

**CLAUDETTE** The convergence plot for the top multi-label runs are in Fig. 31, 32, 33, and 34. For binary, they can be found in in Fig. 35, 36, 37, and 38.

**MRPC** The convergence plot for the top runs are in Fig. 39, 40, 41, and 42.

#### N.3 Diversity Plot

The diversity plot tracks the number of unique in-<br/>dividuals and prompts across generations to assess<br/>genetic diversity. This multi-line plot shows the<br/>unique count of prompt 1, prompt 2, and unique<br/>individuals. The x-axis represents the generation<br/>number, while the y-axis represents the count of<br/>unique individuals. A high count indicates high<br/>diversity, suggesting that the GA is still exploring<br/>solutions, whereas a sharp drop in diversity sug-2064<br/>2065<br/>2065



Figure 9: CLAUDETTE - Left: plot of metrics and average fitness for best multi-label run B in 9. Right: Diversity plotting for best multi-label run B in Table 9



Figure 10: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run B in 9.

gests exploitation, whereby the same individual is being selected for the next generation several times due to high selection pressure. Diversity is crucial for exploration in early generations. The GA may get stuck in a local optimum if diversity drops too early. If diversity remains high for too long, the GA may struggle to converge.

> **CLAUDETTE** Diversity plots for the top multilabel runs are on the right side of Fig. 7, 9 and 11, 13. For the binary runs, diversity plots are on the right of Fig. 15, 17 and 19, 21.

**MRPC** Diversity plots for the top runs are on the right side of Fig. 23, 25 and 27, 29.

#### N.4 Similarity Heatmaps

2074

2075

2079

2084

2090

2091

2094

The similarity heatmap compares the similarity of best individuals across generations using Levenshtein distance. These plots take the form of heatmaps where the x-axis and y-axis represent generations, and the color intensity represents the distance. The darker the color, the more similar (smaller distance) the prompts are. The Levenshtein distance measures character-level differences between best individuals. If distances are high between adjacent generations, it suggests significant mutation and exploration. If distances are low, it suggests convergence and exploitation. Each cell compares the similarity of the best individuals from one generation to another. Diagonal cells should always be darkest since they compare identical generations. Clusters of dark squares suggest stable solution phases in the GA. Although we also plotted the tokenized version of this (where token distance rather than character distance is compared), the plots differ very slightly and globally communicate the same information.

2098

2100

2101

2102

2103

2104

2105

2106

2107

2108

2109

2110

2111

2112

2113

**CLAUDETTE** Prompt similarity plots for the top 4 multi-label runs are in Fig. 8, 10, 12, and 14. For the binary they are in Fig. 16, 18, 20, and 22.

**MRPC** Prompt similarity plots for the top 4 runs are in Fig, 24, 26, 28, and 30.

#### **N.5** Summary of Plot Interpretations

The combination of these plots provides a com-<br/>prehensive view of how the genetic algorithm pro-<br/>gresses over time. The metrics over generations21142115211521162116



Figure 11: CLAUDETTE - Left: plot of metrics and average fitness for best multi-label run C in 9. Right: Diversity plotting for best multi-label run C in 9



Figure 12: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run C in 9.

plot tracks performance trends, the convergence plot highlights stability and volatility, the diversity plot indicates exploration versus exploitation, and the similarity heatmaps reveal how best individuals evolve.

### O Ablation Study

2117

2119

2120

2121

2124

2125

2126

2128

2129

2130

2131

Comparing the pre and post-ablation metric plots (Fig. 43), we observe that the post-ablation plot flatlines for all metrics, including average fitness (and looks similarly flat for the binary case). In contrast, the pre-ablation plot shows a clear trend of exploration and improvement, demonstrating the role of selection in guiding the search toward optimal solutions. By removing it, the evolutionary process collapses into a random stagnating search.

## P LLM-Safe MRPC

2133We performed a thorough preprocessing of the2134Microsoft Research Paraphrase Corpus (MRPC)2135(Dolan and Brockett, 2005) to ensure its suitability2136for modern large language model (LLM) pipelines.2137MRPC consists of sentence pairs extracted from2138news sources, labeled as semantically equivalent

or not. Our preprocessing was carried out with the intent to sanitize potentially problematic content and eliminate parsing issues during downstream processing, which we faced in practice, when we attempted to run our framework on the unprocessed dataset.

2139

2140

2141

2142

2143

2144

2145

2146

2147

2148

2149

2150

2151

2152

2153

2154

2155

2156

2157

2158

#### P.1 Trigger Keyword Removal

We defined a list of content-sensitive trigger keywords that might introduce bias or lead to malformed LLM output due to content flagging. This list included terms such as: ["murder", "terrorist", "rape", "suicide", "nazi", "porn", "overdose", "deep state", ...]

Using a compiled regex, we flagged and removed any sentence pair where either sentence contained one of these keywords. This was applied separately to the training and test sets. We flagged and removed 124 rows from the training set and 53 rows from the test set.

#### P.2 Quote Normalization

Many sentences contained unbalanced or mal-<br/>formed quote characters (e.g., unmatched ", im-<br/>proper smart quotes like " and ", or terminal es-2159<br/>216021612161



Figure 13: CLAUDETTE - Left: plot of metrics and average fitness for best multi-label run D in 9. Right: Diversity plotting for best multi-label run D in 9



Figure 14: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run D in 9.

caped quotes like "). These were identified using a 2162 custom detection function that counted quote occur-2163 2164 rences per sentence and flagged anomalies where the quote count was odd. We manually corrected 2165 374 such cases across both sentence columns. All 2166 forms of quotation marks were then normalized 2167 to a single safe, non-standard Unicode character 2168 (U+2033 Double Prime), visually identical to a 2169 double quote, and interpreted the same by an LLM, 2170 but would not interfere with JSON parsing. 2171

#### P.3 Final Output

2172

2174

2175

2176

2177 2178

- 2173 The final version of the dataset:
  - Contains only rows free of trigger words.
    - Has quote balance issues corrected across all sentence pairs.
  - Is JSON-safe and fully parsable by LLMs and downstream systems.

2179We refer to this cleaned version as the LLM-2180Safe MRPC Dataset and use it consistently2181throughout our experiments.



Figure 15: CLAUDETTE - Left: plot of metrics and average fitness for best binary run A in 8. Right: Diversity plotting for best binary run A in Table 8.



Figure 16: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best binary run A in Table 8.



Figure 17: CLAUDETTE - Left: plot of metrics and average fitness for best run B in 8. Right: Diversity plotting for best binary run B in Table 8.



Figure 18: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best binary run B in Table 8.

Prompt 1	Prompt 2
Create a feature-enriched output that provides a reasoning for each sentence's most likely classification.	For each sentence contained within the input data, evaluate and accurately classify it into one or more of the following categories: ((cat- egory listing)) Carefully analyze the content and implications within each sentence to de- termine the comprehensive set of categories it belongs to.
Generate an explanation-rich classification for each sentence, including the reasoning behind the classification decision.	Analyze each sentence in the input data and classify it into one or more relevant categories based on their content and implications, ensur- ing precision in multi-label classification.
Provide a detailed analysis for each sentence, outlining the reasoning for its classification into the most likely category.	Perform a comprehensive classification of each input sentence into appropriate cate- gories, ensuring all applicable labels are cap- tured.
Construct a comprehensive output that explains the rationale for each sentence's classification.	Evaluate each sentence thoroughly, assigning it to relevant categories and providing precise multi-label classifications.
Develop an enriched response that details the reasoning for each sentence's assigned classification.	Classify the input sentences, ensuring a rigor- ous multi-label classification for relevant as- pects such as: ((category listing))
Offer a feature-oriented output that justifies the classification of each sentence with clear reasoning.	For every sentence in the dataset, determine the applicable categories and provide an accu- rate multi-label classification for these: ((cate- gory listing))
Generate a detailed report justifying each sen- tence's classification with specific reasoning.	Thoroughly analyze each sentence to classify it into one or more relevant categories, captur- ing all dimensions of the classification.
Create a classification output enriched with reasoning for every sentence in the input.	Assign appropriate classifications to each input sentence, reflecting its content and intent while addressing these categories: ((category listing ))
Produce an output that pairs each sentence with an explanation for its classification.	Evaluate and classify each sentence in the dataset into all relevant categories, focusing on ((category listing)).
Develop a thorough output that provides reasoning for the classification of each input sentence.	Analyze the input data sentence by sentence to identify the most applicable categories for each, ensuring completeness in multi-label classification.
Deliver a reasoning-augmented classification output for each provided sentence.	Classify the content of each sentence with a focus on accurate multi-label categorization, rigorously addressing ((category listing)).

Table 4: CLAUDETTE - Manual multi-label prompt bank used to initialize every GenDLN multi-label run.

Prompt 1	Prompt 2
You are a linguistic analysis model special-	You are an expert in paraphrase detection. In
ized in paraphrase tasks. For each input pair,	the following your task is to analyze if sen-
extract key semantic and syntactic features rel-	tence 2 is a paraphrased version of sentence
evant for paraphrase classification.	1. Thus, you shall classify each sentence pair
	into 0 ('not equivalent') or 1 ('equivalent')
	depending on whether sentence 1 and 2 are
	semantically equivalent.
Analyze each sentence pair to identify mean-	Given each sentence pair, determine if the sec-
ingful features that help determine if the two	ond sentence is a paraphrase of the first. Out-
sentences are paraphrases.	put 1 if they are semantically equivalent, 0 if
	they are not.
Given a list of sentence pairs, extract discrimi-	Your job is to judge whether the meaning of
native features for each pair that can support	sentence 1 is preserved in sentence 2. Clas-
downstream paraphrase detection.	sify the pair as 1 for paraphrase or 0 for non-
	paraphrase.
You are tasked with analyzing sentence pairs.	Classify each sentence pair by checking if sen-
For each pair, return a compact description of	tence 2 can be considered a paraphrase of sen-
important features that would help in classify-	tence 1. Use 1 for equivalent, 0 for not equiva-
ing paraphrase relationships.	lent.
Analyze the input sentence pairs and extract	You are a paraphrase classification assistant.
useful features that would support a classifier	For each sentence pair, assign a binary label:
in detecting semantic equivalence.	1 if sentence 2 is a paraphrase of sentence 1,
	else 0.
You are a feature extraction system for para-	You are to detect paraphrases. For each sen-
phrase detection. For each sentence pair, out-	tence pair, determine if both express the same
put key comparison features in the specified	meaning. Label with 1 if equivalent, otherwise
format.	0.
Given sentence pairs, identify and summarize	For each given pair of sentences, assess
linguistic or semantic cues that are relevant for	whether sentence 2 paraphrases sentence 1.
determining paraphrasing.	Output 1 for equivalent meaning, 0 for dif-
	ferent meaning.
For each pair of sentences, write a brief set	You are evaluating sentence-level semantic
of features that capture their semantic, lexical,	similarity. Classify each pair with 1 if both
and structural alignment.	sentences are paraphrases, and 0 if they are
	not.
Table 5: MRPC - Manual binary prompt bank ()	Part 1/3) used to initialize GenDLN binary runs.



Figure 19: CLAUDETTE - Left: plot of metrics and average fitness for best run C in 8. Right: Diversity plotting for best binary run C in Table 8.

Prompt 1	Prompt 2
Inspect each input sentence pair and generate	You are an NLP expert assessing paraphrase
their similarity or difference in meaning.	semantically equivalent, else 0.
You are a natural language understanding	You are a binary classifier for sentence equiv-
model. For each sentence pair, extract features	alence. Judge whether sentence 2 retains the
that reveal differences or overlaps in meaning and expression.	meaning of sentence 1. Output 1 or 0 accordingly.
Identify semantic relationships and stylistic	Your goal is to assess if sentence 2 can be con-
variations in each sentence pair. Output con-	sidered a reasonable paraphrase of sentence 1.
cise features that explain their alignment or divergence.	Output 1 if so, otherwise 0.
For every input pair, generate a feature-based	Examine the semantic content of each sentence
comparison that highlights differences in struc-	pair and decide if they convey the same core
ture, meaning, or terminology.	meaning. Return 1 for paraphrase, 0 for other- wise.
You are helping a classifier understand sen-	Determine whether sentence 2 is interchange-
tence similarity. Extract key features that	able with sentence 1, i.e. a suitable paraphrase.
could guide a model in deciding paraphrase	Output 1 if they are interchangeable, else 0.
equivalence.	Vou are accessing perperture validity. Classify
ings nuanced differences or structural shifts	each pair as 1 if the second sentence accurately
Provide these insights as short structured fea-	reflects the meaning of the first or 0 if not
tures.	reneets the meaning of the mist, of o it not
Your goal is to support a paraphrase detec-	For every pair, identify whether sentence 2 ex-
tion system by extracting features that capture	presses the same meaning as sentence 1 using
lexical, syntactic, and semantic properties of sentence pairs.	a binary label: 1 (yes), 0 (no).
Review each sentence pair and write a concise	Your task is to judge if sentence 2 carries the
summary of alignment cues and linguistic dif-	same intent and meaning as sentence 1. Output
ferences that may affect paraphrase detection.	1 for equivalence, 0 otherwise.

Table 6: MRPC - Manual binary prompt bank (Part 2/3) used to initialize GenDLN binary runs.



Figure 20: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best binary run C in Table 8.

Prompt 1	Prompt 2
As a sentence-level feature extractor, outline	Determine semantic equivalence at the sen-
the textual signals that could be used to deter-	tence level. For each pair, output 1 if meaning
mine if two statements express the same idea.	is preserved between the two sentences, 0 if it is lost or altered.
Examine each sentence pair and extract dis-	Review each sentence pair and determine
tinguishing features that would help a down- stream model judge paraphrase likelihood.	whether sentence 2 retains the essential mean- ing of sentence 1. Respond with 1 for equiva- lence, 0 otherwise.
Your job is to find patterns in sentence pairs	Your job is to classify whether sentence 2 can
that indicate whether they express similar or	logically be interpreted as expressing the same
different meanings. Output a compact list of relevant features.	idea as sentence 1. Output 1 for yes, 0 for no.
You are a linguistic alignment engine. Identify whether key predicates, named entities, and relationships are preserved across the sentence pair.	Assess whether sentence 2 paraphrases sen- tence 1 without introducing or omitting criti- cal information. Output 1 for paraphrase, 0 if meaning changes.
Highlight phrasing shifts, information asym- metry, or reordering patterns that could influ- ence whether the sentence pair is semantically aligned.	For each pair of statements, decide whether sentence 2 communicates the same content as sentence 1. Respond with 1 for equivalent, 0 for not equivalent.
For each input pair, extract lexical and struc-	Analyze the sentence pair and determine if
tural markers - including synonym usage, clause structure, and entity alignment - that contribute to paraphrase detection.	their meanings align well enough to be consid- ered paraphrases. Output 1 if they do, 0 if not.
Extract the central premise of each of the two sentences, what information does each convey?	Are they paraphrases of each other? Output 1 for yes, 0 for no.
As an expert writer, would you say the two sentences convey the same main idea? What would you say is the point of each sentence?	Would it be reasonable to replace one sentence with the other in a text without changing the overall meaning? In other words, are the sen- tences paraphrases of each other? Output 1 if yes and 0 if no.
Could the two sentence reasonably be ex- changed within a text without changing the general meaning of the text? Why or why not?	Given that assessment, can the sentences be classified as paraphrases of each other? An- swer with 1 if they are paraphrases, and 0 if not.
Table 7: MRPC - Manual binary prompt bank (         Best Individual Metrics by Generation vs. Average Fitness	Part 3/3) used to initialize GenDLN binary runs.



Figure 21: CLAUDETTE - **Left:** plot of metrics and average fitness for best run D in 8. **Right:** Diversity plotting for best binary run D in Table 8.

Metric	Run A	Run B	Run C	Run D
Runtime (mins)	58.565	160.9097	100.8069	53.262
Best Fitness	0.8785	0.8785	0.8687	0.8380
Best Accuracy	0.8788	0.8788	0.8687	0.8384
Test. Accuracy	0.7897	0.7706	0.7646	0.7404
Best Macro F1	0.8785	0.8785	0.8686	0.8380
Test. Macro F1	0.6523	0.6364	0.6338	0.6172
Best Weighted F1	0.8784	0.8784	0.8687	0.8379
Test. Weighted F1	0.8256	0.8115	0.8073	0.7894
Selection Strategy	Rank	SUS	SUS	Rank
Crossover Type	Semantic Blending	Token Level	Semantic Blending	Semantic Blending
Crossover Rate	0.800	0.800	0.800	0.800
Mutation Type	Semantic	Syntactic	Semantic	Semantic
Mutation Rate	0.200	0.200	0.200	0.200
Population Size	10	30	30	10
<b>Completed Generations</b>	16	16	9	16
Stopped Early	Yes	Yes	Yes	Yes
Stopped Early Reason	5 stag. gens.	5 stag. gens.	5 stag. gens.	5 stag. gens.

Table 8: CLAUDETTE - Selected runs for binary (fair/unfair) classification.

Metric	Run A	Run B	Run C	Run D
Runtime (mins)	469.689	439.694	373.876	155.367
Best Fitness	0.938	0.925	0.922	0.921
Best Accuracy	0.910	0.890	0.880	0.900
Test. Accuracy	0.825	0.769	0.809	0.802
Best Macro F1	0.947	0.936	0.935	0.929
Test. Macro F1	0.862	0.799	0.844	0.855
Best Weighted F1	0.944	0.933	0.929	0.923
Test. Weighted F1	0.856	0.808	0.842	0.851
Selection Strategy	Rank	Steady-State	SUS	Steady-State
Crossover Type	Phrase Swap	Phrase Swap	Token Level	Semantic Blending
Crossover Rate	0.850	0.850	0.850	0.800
Mutation Type	Insertion	Insertion	Syntactic	Semantic
Mutation Rate	0.300	0.300	0.300	0.200
Population Size	30	30	30	30
<b>Completed Generations</b>	30	30	30	12
Stopped Early	No	No	No	Yes
Stopped Early Reason	-	-	-	5 stag. gens.

Table 9: CLAUDETTE - Selected best runs for multi-label classification.

Metric	Run A	Run B	Run C	Run D
Runtime (mins)	137.681	167.228	67.070	127.262
Best Fitness	0.850	0.840	0.850	0.840
Best Accuracy	0.850	0.840	0.850	0.840
Test. Accuracy	0.813	0.807	0.798	0.799
Best Macro F1	0.850	0.840	0.850	0.840
Test. Macro F1	0.796	0.787	0.782	0.781
Best Weighted F1	0.850	0.840	0.850	0.840
Test. Weighted F1	0.816	0.809	0.802	0.802
Selection Strategy	Steady-State	Roulette	Tournament	SUS
Crossover Type	Single Point	Semantic Blending	Token Level	Two Point
Crossover Rate	0.85	0.85	0.85	0.80
Mutation Type	Semantic	Insertion	Insertion	Deletion
Mutation Rate	0.20	0.20	0.20	0.20
Population Size	30	30	30	30
<b>Completed Generations</b>	16	23	12	15
Stopped Early	Yes	Yes	Yes	Yes
Stopped Early Reason	10 stag. gens.	10 stag. gens.	10 stag. gens.	10 stag. gens.

Table 10: MRPC - Selected best runs for binary paraphrase classification.

Run	Prompt Text
А	<b>Prompt 1</b> : Assess the potential legal consequences and issues of the following sentence.
	<b>Prompt 2</b> : Based on the previous discussion, would you consider this sentence to be fair or
	unfair as it stands?
В	<b>Prompt 1</b> : Interpret the following sentence in any hidden clauses or implications.
	<b>Prompt 2</b> : Will the described potential impact be considered fair or unfair?
С	Prompt 1: Assess the possible legal ramifications and effect on consumer rights of the
	following sentence.
	<b>Prompt 2</b> : Considering the impact of the ethical implications discussed, is this sentence fair
	or unfair in its current phrasing?
D	<b>Prompt 1</b> : Identify any potential legal issues when analyzing the meaning of the following
	sentence in a legal context.
	<b>Prompt 2</b> : Given the emphasized issues, is this sentence fair or unfair in its current state?

Table 11: CLAUDETTE - Prompt 1 and 2 of the best individuals for the runs as reported in Table 8 for the binary classification task.



Figure 22: CLAUDETTE - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best binary run D in Table 8.

Run	Prompt Text
А	<b>Prompt 1</b> : To enhance transparency for the end user, who may not be familiar with the
	internal mechanics of our system, we should annotate each individual sentence contained
	within the given customer review that is specifically about our recently introduced product,
	including a clear, concise, and straightforward explanation that meticulously details the
	reasoning, justification, and rationale behind its specific classification, ensuring that the user
	comprehends why we classified the sentence as such.
	<b>Prompt 2</b> : To thoroughly organize and accurately assign a precise data monitoring technique
	or pertinent cookie policies that are explicitly outlined in a legal privacy policy document, a
	team of legal experts should meticulously review the entire policy document, starting from
	the introduction to the conclusion, and systematically classify each individual clause from the
	contract with high precision during the detailed multi-label classification process, ensuring
	that the resulting labels are not only relevant to the contractual obligations clearly outlined in
	the legal documents but also precise in their legal definition.
В	<b>Prompt 1</b> : To ensure thorough documentation and transparency in our contractual legal anal-
	ysis efforts within the jurisdiction of the relevant state legal system, produce a comprehensive
	legal classification of the content within each individual clause that is clearly outlined in the
	case files pertaining to the ongoing corporate lawsuit.
	<b>Prompt 2</b> : When examining corporate legal documents, such as those related to IT service
	agreements, systematically classify each individual sentence from various types of contrac-
	tual clauses, including confidentiality, liability, and termination clauses, into relevant and
	predefined labels for better organization and analysis.
С	<b>Prompt 1</b> : Present a detailed report on the categorization of every sentence, accompanied by
	relevant evidence.
	<b>Prompt 2</b> : Every sentence, in the multi-label classification process, will be assigned to its
	fitting categories to maintain it thoroughly, emphasizing suitable labels that range from PINC
	for cookie and tracking to LAW for legal frameworks.
D	<b>Prompt 1</b> : Generate a feature-focused output that matches each sentence with a reason for
	its categorization.
	<b>Prompt 2</b> : Sort and classify each sentence in the dataset, taking into account these categories:
	PINC (Cookies or data collection), USE (Rules on user activities), CR (Removal rights),
	TER (Service terminations), LTD (Limitation of liability), A (Arbitration resolutions), LAW
	(Governing legal codes), J (Jurisdiction clauses), CH (Agreement changes).

Table 12: CLAUDETTE - Prompt 1 and 2 of the best individuals for the runs as reported in Table 9 for the multi-label task.



Figure 23: MRPC - Left: plot of metrics and average fitness for best run A in Table 10. Right: Diversity plotting for best multi-label run A in Table 10

Run	Prompt Text
А	<b>Prompt 1</b> : Assess each pair of sentences and generate a feature-based comparison that
	highlights differences in structure, meaning, or terminology.
	<b>Prompt 2</b> : You are evaluating each pair of sentences to determine if they express the same
	central meaning; return 1 if they are paraphrases, and 0 otherwise.
В	<b>Prompt 1</b> : For each individual pair of sentences that you evaluate within a comparative text
	analysis study, output a meaningful feature description that accurately captures their shared
	meanings, specific word choices, sentence structure, and stylistic differences.
	<b>Prompt 2</b> : After carefully examining each individual pair of sentences for their meaning
	and content, determine if they are paraphrases and convey the same meaning; label with a 1
	if they are semantically equivalent, otherwise label them with a 0.
С	<b>Prompt 1</b> : For each sentence pair, extract semantic relationships and output concise features
	that reveal differences or overlaps in meaning and expression.
	<b>Prompt 2</b> : Your goal is to assess whether or not sentence 2 retains the meaning of sentence
	1, taking into account all aspects of semantics and context. Judge whether sentence 2 can
	be considered a reasonable paraphrase of sentence 1, with an equivalent core interpretation.
	Output 1 for yes or 0 for no accordingly.
D	<b>Prompt 1</b> : Compare each sentence pairs that reveal distinguishing features in meaning.
	<b>Prompt 2</b> : Judge whether they are expressing the same intent of each other in a text.

Table 13: MRPC - Prompt 1 and 2 of the best individuals for the runs as reported in Table 10 for the paraphrase classification task.



Figure 24: MRPC - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run A in Table 10.



Figure 25: MRPC - Left: plot of metrics and average fitness for best run A in Table 10. Right: Diversity plotting for best multi-label run B in Table 10



Figure 26: MRPC - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run B in Table 10.



Figure 27: MRPC - Left: plot of metrics and average fitness for best run C in Table 10. Right: Diversity plotting for best multi-label run C in Table 10



Figure 28: MRPC - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run C in Table 10.



Figure 29: MRPC - Left: plot of metrics and average fitness for best run D in Table 10. Right: Diversity plotting for best multi-label run D in Table 10



Figure 30: MRPC - Best Prompt 1 and Prompt 2 Levenshtein distance matrix across generations for best multi-label run D in Table 10.



Fitness Convergence Over Generations 0.90 80.85 Eitness 08.0 0.75 10 Generation est Fitness · · Worst Fitness ga\_log\_20250120\_160348.log

Figure 31: CLAUDETTE - Convergence plot for best multi-label run A in Table 9. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.

Figure 34: CLAUDETTE - Convergence plot for best multi-label run D in Table 9. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.



Figure 32: CLAUDETTE - Convergence plot for best multi-label run B in Table 9. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.



Figure 33: CLAUDETTE - Convergence plot for best multi-label run C in Table 9. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.



Figure 35: CLAUDETTE - Convergence plot for best binary run A in Table 8.



Figure 36: CLAUDETTE - Convergence plot for best binary run B in Table 8.



Fitness Convergence Over Generations

Figure 37: CLAUDETTE - Convergence plot for best binary run C in Table 8. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.

Figure 40: MRPC - Convergence plot for best binary run B in Table 10. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.



Figure 38: CLAUDETTE - Convergence plot for best binary run D in Table 8.



Figure 39: MRPC - Convergence plot for best binary run A in Table 10. X on the x-axis indicates an illegal individual as per the fallback mechanism in Appendix D.



Figure 41: MRPC - Convergence plot for best binary run C in Table 10.



Figure 42: MRPC - Convergence plot for best binary run D in Table 10.



Figure 43: CLAUDETTE - Left: plot of metrics and average fitness for best multi-label run in Table 1. Right: Ablation of selection pressure for the same run.