

Enhancing LLM-Based Neural Network Generation: Few-Shot Prompting and Efficient Validation for Automated Architecture Design

Anonymous CVPR submission

Paper ID

Abstract

001 While large language models (LLMs) have unlocked a
002 new paradigm for automated neural architecture design,
003 their generative power remains poorly understood: the crit-
004 ical question of how in-context example count governs gen-
005 eration quality has never been systematically investigated.
006 Building on the NNGPT/LEMUR framework, we conduct
007 the first empirical study of few-shot prompting for neural
008 architecture generation. We generate and evaluate 1,900
009 neural architectures produced by LLMs across six com-
010 mon vision benchmarks, rigorously assessing how the num-
011 ber of supporting examples ($n \in \{1, \dots, 6\}$) shapes gen-
012 eration stability, architectural diversity, and early-epoch
013 validation performance. Although $n=3$ yields the high-
014 est dataset-balanced mean accuracy (53.1%), this improve-
015 ment is task-dependent. Performance gains increase with
016 task complexity (CIFAR-100: +11.6%, $p=0.001$, $d=0.73$;
017 CelebA-Gender: +6.5% at $n=2$, $p=0.038$), whereas larger
018 context sizes ($n>3$) lead to statistically significant degrada-
019 tion on more structured benchmarks (ImageNette: -14.5%,
020 $p=0.010$; CIFAR-10: -8.4%, $p=0.016$). At $n=6$, gener-
021 ation performance collapses entirely (99.8% failure rate),
022 further corroborating that $n=3$ strikes an optimal bal-
023 ance between providing informative context and avoiding
024 prompt saturation. Qualitative analysis reveals that $n=3$
025 uniquely enables architectural pattern synthesis—hybrid
026 ResNet-DPN and ResNet-AlexNet structures—absent un-
027 der single-example prompting. To support scalable de-
028 ployment, we introduce whitespace-normalized hashing for
029 real-time duplicate detection, achieving a $100\times$ speedup
030 over AST-based methods and eliminating redundant train-
031 ing of formatting-level duplicates. Together, these findings
032 establish prompt context calibration and lightweight valida-
033 tion as key principles for scalable, resource-efficient LLM-
034 driven architecture design.

1. Introduction

035
036 The design of effective neural network architectures tradi-
037 tionally relies on domain expertise and iterative experimen-
038 tation. Neural Architecture Search (NAS) [26, 33] methods
039 automate this process, but often require substantial compu-
040 tational resources and carefully engineered search spaces.
041 Recently, Large Language Models (LLMs) have emerged
042 as a promising alternative paradigm, leveraging code gener-
043 ation capabilities to synthesize neural architectures directly
044 from natural language prompts [5, 30]. This approach shifts
045 architecture design from gradient-based search toward gen-
046 erative modeling of network structure.

047 **The Promise and Gap of LLM-Based Architecture**
048 **Search.** Recent advances have demonstrated that LLMs
049 can serve as architecture generators through code synthe-
050 sis. EvoPrompting [3] combines evolutionary algorithms
051 with prompt-tuned LLMs, achieving competitive results on
052 MNIST-1D and CLRS benchmarks. LLMatic [22] applies
053 quality-diversity search with LLMs to generate diverse ar-
054 chitectures on CIFAR-10. The NNGPT framework [12]
055 demonstrates end-to-end LLM-driven architecture synthe-
056 sis using the LEMUR dataset [9].

057 However, these works adopt different prompting strate-
058 gies without systematic investigation of a fundamental pa-
059 rameter: *how many in-context examples optimize genera-
060 tion quality?* EvoPrompting uses evolutionary selection of
061 examples but does not isolate the effect of example quantity.
062 NNGPT employs single-reference prompting. LLMatic fo-
063 cuses on quality-diversity optimization rather than prompt
064 engineering. This gap is critical because few-shot learning
065 is known to be highly sensitive to example count in natu-
066 ral language tasks [19, 21], yet neural architecture genera-
067 tion presents unique constraints: architectures are substan-
068 tially longer than typical code snippets, architectural diver-
069 sity matters alongside functional correctness, and context
070 window limitations become critical.

071 **Research Objective.** This work provides the first sys-
072 tematic empirical study of few-shot example count in LLM-
073 based neural architecture generation. We address three

074 questions:

075 **RQ1 (Context Capacity):** How does the number of in-
076 context architectural examples influence generation stabil-
077 ity and early-epoch performance?

078 **RQ2 (Architectural Synthesis):** Does moderate con-
079 text enrichment enable qualitatively different architectural
080 patterns compared to single-example prompting?

081 **RQ3 (Computational Efficiency):** Can lightweight
082 deduplication strategies reduce redundant training in large-
083 scale LLM-based generation pipelines?

084 To investigate these questions, we build upon the NNGP-
085 T/LEMUR [12] framework and conduct a controlled study
086 varying the number of supporting architectures provided in
087 the prompt ($n \in \{1, \dots, 6\}$). Across seven computer vision
088 benchmarks, we generate and train 1,900 unique architec-
089 tures under a rapid screening protocol. Our results reveal
090 a stable regime at three in-context examples and a marked
091 decline in generation success for larger context sizes, indi-
092 cating practical limits on prompt enrichment.

093 In addition, large-scale generation pipelines introduce a
094 practical challenge: LLMs frequently produce formatting-
095 level duplicates that lead to redundant training. To mitigate
096 this issue, we adopt a lightweight syntactic canonicaliza-
097 tion strategy based on whitespace normalization and hash-
098 ing, enabling efficient detection of duplicate code prior to
099 training. While this approach does not capture semantic
100 equivalence, it eliminates formatting-based redundancy at
101 negligible computational cost.

102 **Key Findings.** Our large-scale experiments across 1,900
103 architectures and seven benchmarks reveal four principal re-
104 sults. Three supporting examples ($n=3$) achieves the high-
105 est dataset-balanced mean accuracy (53.1%), with statisti-
106 cally significant gains on fine-grained tasks (CIFAR-100:
107 +11.6%, $p=0.001$, Cohen’s $d=0.73$), while generation suc-
108 cess collapses at $n=6$ with a 99.8% failure rate — estab-
109 lishing a hard empirical upper bound on prompt context
110 size. Qualitative analysis shows that $n=3$ uniquely en-
111 ables cross-paradigm synthesis: hybrid ResNet-DPN and
112 ResNet-AlexNet architectures emerge that are entirely ab-
113 sent under single-example prompting, confirming that mod-
114 erate context enrichment unlocks qualitatively richer design
115 space exploration. Finally, whitespace-normalized hashing
116 achieves a $100\times$ speedup over AST-based deduplication at
117 <1 ms per sample, enabling real-time duplicate rejection
118 in any large-scale LLM-based generation pipeline. Over-
119 all, these findings establish prompt context calibration and
120 lightweight validation as two key principles for scalable,
121 resource-efficient LLM-driven architecture design, comple-
122 menting existing approaches such as EvoPrompting [3] and
123 LLMatic [22] that focus on search strategy rather than
124 prompt engineering.

2. Related Work 125

2.1. Neural Architecture Search 126

Traditional NAS methods employ reinforcement learn- 127
ing [33], evolutionary algorithms [26], or gradient- 128
based optimization [17] to discover effective architectures. 129
ENAS [25] improves efficiency through weight sharing, re- 130
ducing search costs from thousands to single-digit GPU 131
days. However, these approaches require carefully designed 132
search spaces and substantial computational resources. 133

2.2. LLMs for Code Generation 134

Foundation models trained on code repositories [4, 15, 27] 135
have demonstrated remarkable code generation capabili- 136
ties. DeepSeek Coder [7], trained on 2 trillion tokens from 137
87 programming languages, achieves state-of-the-art per- 138
formance on code generation benchmarks. CodeGen [24] 139
demonstrates that domain-specific fine-tuning on scientific 140
computing libraries (PyTorch, NumPy) further improves 141
generation quality for technical domains. 142

2.3. LLM-Based Architecture Generation 143

Recent work explores using LLMs as architecture gener- 144
ators, offering computational advantages over traditional 145
NAS. 146

EvoPrompting [3] combines evolutionary algorithms 147
with LLM code generation for architecture search. Using 148
soft prompt-tuning and iterative refinement across gener- 149
ations, EvoPrompting outperforms human-designed base- 150
lines on MNIST-1D and CLRS benchmarks. However, their 151
few-shot prompting strategy uses evolutionary population- 152
based selection without systematic investigation of optimal 153
example count. 154

LLMatic [22] integrates LLMs with MAP-Elites 155
quality-diversity optimization for architecture search on 156
CIFAR-10 and NAS-Bench-201. Evaluating $\sim 2,000$ archi- 157
tectures, their approach emphasizes architectural diversity 158
through multi-objective optimization but does not investi- 159
gate prompt engineering strategies. 160

NNGPT [12] introduces an end-to-end framework for 161
LLM-driven neural architecture synthesis, built on the 162
LEMUR dataset [9]. NNGPT demonstrates one-shot gener- 163
ation of complete training specifications including archi- 164
tecture and hyperparameters. While powerful, NNGPT’s 165
prompting strategy relies on single reference models, leav- 166
ing unexplored the impact of multiple supporting examples 167
on generation quality and architectural diversity. 168

Positioning Our Work. While these approaches demon- 169
strate LLM feasibility for architecture generation, *none* 170
systematically investigate the impact of few-shot example 171
count. EvoPrompting uses evolutionary selection of exam- 172
ples but does not isolate the effect of example quantity. 173
NNGPT uses single-reference prompting ($n = 1$). LLMatic 174

175 focuses on quality-diversity optimization rather than prompt
176 configuration. Our work fills this gap through controlled ex-
177 periments varying $n \in \{1, 2, 3, 4, 5, 6\}$, identifying $n = 3$
178 as optimal for balancing performance, diversity, and gener-
179 ation stability in computer vision tasks.

180 2.4. Few-Shot Prompting Strategies

181 Few-shot in-context learning [2] has become fundamental
182 to LLM applications. Recent work reveals that demon-
183 stration count significantly impacts performance: Min et
184 al. [21] find that 4-8 examples optimize most tasks, while
185 Lu et al. [19] show that example ordering matters critically.
186 Chain-of-thought prompting [30] further enhances reason-
187 ing capabilities through structured demonstrations.

188 However, these findings focus on natural language and
189 general programming tasks. Neural architecture generation
190 presents unique constraints: (1) architectures are substan-
191 tially longer (~ 500 tokens) than typical code snippets (~ 50
192 tokens), increasing context pressure; (2) architectural diver-
193 sity matters alongside functional correctness; (3) context
194 window limitations become critical. Our work establishes
195 that $n = 3$ balances these competing demands for vision ar-
196 chitectures, extending few-shot prompting theory to neural
197 synthesis.

198 2.5. Code Deduplication and Similarity Detection

199 Traditional approaches use AST-based comparison [28] or
200 graph-based representations [10], providing high accuracy
201 but incurring 10-100ms overhead per comparison. While
202 Winnowing [29] offers faster fingerprinting through token-
203 based hashing, it misses formatting variations common in
204 LLM-generated code.

205 Our whitespace-normalized hashing targets the domi-
206 nant source of duplication in LLM pipelines—formatting
207 variations—achieving < 1 ms computation while maintain-
208 ing zero false negatives on 1,900 generated architec-
209 tures. While more sophisticated semantic equivalence de-
210 tection [10] could complement our approach, the $100\times$
211 speedup vs. AST parsing makes syntactic canonicalization
212 practical for real-time generation workflows.

213 3. Methodology

214 3.1. NNGPT Base System

215 Inspired by recent advancements in the application of
216 LLMs across various domains, and leveraging the exist-
217 ing LEMUR dataset ecosystem [9], we developed our solu-
218 tion upon the NNGPT architecture generation pipeline [12],
219 which consists of:

220 **Base Language Model:** DeepSeek Coder 7B [7], a
221 state-of-the-art code generation model trained on 2 trillion
222 tokens with specialized focus on Python scientific comput-
223 ing libraries.

LoRA Fine-Tuning: Low-Rank Adaptation [11] with
rank $r = 32$, $\alpha=32$, enabling efficient fine-tuning on
the LEMUR dataset with only ~ 35 M trainable parameters
($\sim 0.5\%$ of base model).

Generation Parameters: Temperature=0.6, top-k=50,
top-p=0.95, max_tokens=65,536, balancing creativity and
syntactic correctness.

3.2. Few-Shot Architecture Prompting (FSAP)

Building upon recent LLM applications in computer vi-
sion and leveraging the LEMUR ecosystem [9], we adopt a
structured prompt construction strategy to isolate the effect
of in-context example count on neural architecture synthe-
sis.

Unlike prior work that employs evolutionary selec-
tion [3] or quality-diversity optimization [22], our goal is
to systematically characterize how context size alone influ-
ences generation quality. By controlling all other variables
(base LLM, fine-tuning strategy, generation temperature),
we provide the first empirical analysis of few-shot scaling
effects in architecture generation.

3.2.1 Prompt Construction

Our prompting strategy is designed to isolate context size
as the sole independent variable. Given a target dataset \mathcal{D}
and desired example count $n \in \{1, \dots, 6\}$, we query the
LEMUR database for top-performing architectures on \mathcal{D} ,
select one as the reference model, and sample the n support-
ing examples uniformly at random. This random sampling
deliberately avoids performance-based heuristics, ensuring
observed differences reflect context quantity rather than ex-
ample quality.

The prompt combines four components: (1) task descrip-
tion and dataset specifications, (2) the reference architecture
with complete code and reported accuracy, (3) n supporting
architectures with code and performance metrics, and (4)
explicit generation constraints.

This controlled design contrasts with prior frameworks:
unlike EvoPrompting [3], which conflates quantity and
quality signals through evolutionary selection, we hold
quality fixed; unlike NNGPT [12], which uses a fixed
prompt ($n=1$), we systematically sweep $n \in \{1, \dots, 6\}$;
and unlike LLMatic [22], which focuses on quality-
diversity search algorithms, we treat the prompt itself as the
primary experimental variable. This provides insights ap-
plicable to any LLM-based generation pipeline regardless
of search strategy. By controlling all other variables base
LLM, LoRA fine-tuning configuration, generation temper-
ature, and dataset split we provide what is, to our knowl-
edge, the first empirical analysis of few-shot scaling ef-
fects specifically in neural architecture generation, a domain
where the unique constraints of long code outputs (~ 500

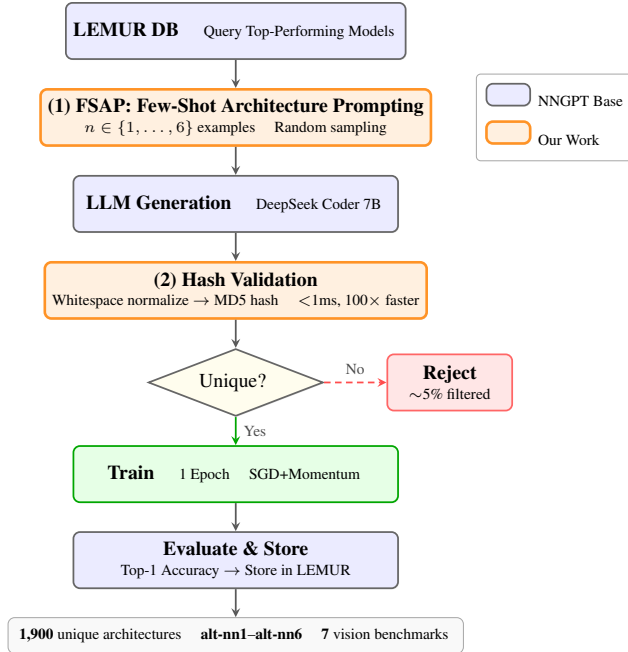


Figure 1. Architecture generation pipeline integrating our contributions into NNGPT. **Orange boxes** highlight our work: (1) FSAP systematically varies supporting example count $n \in \{1, \dots, 6\}$; (2) Whitespace-normalized hash validation achieves $100\times$ speedup over AST-based deduplication. The pipeline processes 1,900 unique architectures across six benchmarks, rejecting $\sim 5\%$ formatting duplicates before training.

tokens per architecture) and strict syntactic validity requirements make context management qualitatively different from standard natural language few-shot settings [19, 21].

3.3. Whitespace-Normalized Hash Validation

Large-scale LLM generation frequently produces formatting-level duplicates that cause redundant training. Our deduplication pipeline applies three steps: whitespace removal to create a formatting-invariant representation, MD5 hashing to produce a unique identifier, and B-tree indexed LEMUR lookup to check existence. The procedure runs in $O(|\mathcal{C}| + \log N)$ time, achieving < 1 ms per sample and a $100\times$ speedup over AST-based alternatives, with zero false positives across 4,033 generated architectures. Full pseudocode and timing analysis are provided in the supplementary material (Algorithm S2, Table S1).

3.4. Integrated Pipeline

Figure 1 illustrates the complete architecture generation pipeline with our contributions integrated into NNGPT.

The pipeline operates as follows:

1. **FSAP** constructs enriched prompts with n supporting examples
2. **LLM Generation** produces candidate architecture code

3. **Hash Validation** checks uniqueness before training
4. **Rejected duplicates** save 2 to 3 GPU hours per instance
5. **Accepted architectures** proceed to training and evaluation
6. **Validated models** are stored in LEMUR with hash identifiers

4. Experiments

4.1. Experimental Setup

Architecture Variants All generated architectures use the prefix `alt-nn`. Variants `alt-nn1` through `alt-nn6` correspond to prompts containing $n = 1$ through $n = 6$ supporting examples, respectively. This controlled variation enables analysis of context size effects under otherwise identical generation settings.

Datasets We evaluate across six computer vision benchmarks of varying scale and difficulty:

- MNIST [14] (10 classes)
- CelebA-Gender [18] (binary classification)
- CIFAR-10 [13] (10 classes)
- CIFAR-100 [13] (100 classes)
- ImageNette (10-class ImageNet subset)
- SVHN [23] (10 classes)

These datasets span simple digit recognition to challenging classification, allowing evaluation under heterogeneous task complexity. **Note:** A seventh benchmark, Places365 [32] (365 classes), was excluded from quantitative comparison tables due to substantially longer training times (180 minutes per epoch), but was used to validate hash-based deduplication efficiency at scale.

Rapid Screening Protocol Our objective is to analyse relative trends across prompt variants at large scale (1,900 architectures). To enable broad comparative evaluation under constrained compute budgets, we adopt a rapid screening protocol consisting of a single training epoch per model.

This setting provides an early-performance signal [8] that allows controlled comparison of context-size variants, rather than serving as an estimate of final convergence performance. All variants are evaluated under identical optimization settings.

Optimization Details

To evaluate our approach rigorously, we adopt true validation accuracy after the first training epoch as our primary metric [8], rather than relying on zero-shot NAS proxies, which, although correlated with fully trained accuracy [6, 16, 20], remain indirect indicators of performance. Even the strongest reported Spearman rank-correlation coefficients ($\rho \approx 0.5\text{--}0.82$) on standard benchmarks such as NAS-Bench-101 and NAS-Bench-201 correspond to a coefficient of determination of at most $R^2 \approx 0.67$, leaving substantial variance unexplained [1, 31]. By using first-epoch validation accuracy, we aim to demonstrate that our

algorithm can reliably influence and accelerate early-stage performance trajectories of neural networks.

- Epochs: 1
- Optimizer: SGD with momentum
- Batch size: dataset-dependent (64-4096)
- Metric: Top-1 accuracy

Generation Statistics

- Total unique architectures trained: 1,900
- Architecture variants: alt-nn1 through alt-nn6 ($n = 1$ to $n = 6$)
- Data transformation variants: 120
- Hash-based validation: Applied to all generated code prior to training

The hash-based duplicate filtering prevents redundant training of formatting-level duplicates (identical architectures with different whitespace/indentation), ensuring computational efficiency across all dataset scales.

4.2. Evaluation Metrics

Primary Metric: Top-1 accuracy after 1 epoch

Statistical Validation: Independent t-tests comparing variants within each dataset, dataset-balanced means for overall comparison

Efficiency Metrics: Hash computation time, duplicates detected

While single epoch evaluation is inherently noisy, our analysis focuses on consistent trends across prompt variants rather than precise architecture ranking.

5. Results and Discussion

5.1. Dataset-Balanced Evaluation

Because the number of successfully generated architectures varied across datasets and prompt variants, naive aggregation of all samples could bias results toward easier tasks (e.g., MNIST) and underrepresent more challenging datasets (e.g., CIFAR-100 or Places365).

To mitigate this imbalance, we adopt a dataset-balanced aggregation strategy. Specifically, we first compute the mean accuracy for each (variant, dataset) pair and then average these per-dataset means, assigning equal weight to each benchmark. This corresponds to macro-averaging across datasets.

Statistical significance is assessed within each dataset using independent t-tests ($p < 0.05$), ensuring that comparisons are not confounded by cross-dataset difficulty differences.

5.2. Few-Shot Prompting Performance

Figure 2 illustrates the relationship between supporting example count and both performance and generation success.

Table 1 presents overall performance across all datasets using balanced statistics.

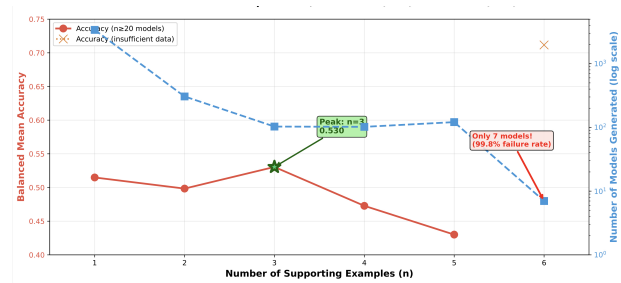


Figure 2. Effect of context size on performance and generation success. Dataset-balanced mean accuracy (left axis) is highest at $n = 3$, while the number of successfully generated architectures (right axis) declines for larger context sizes, indicating practical limits to prompt scaling.

Table 1. Dataset-balanced macro-mean across six benchmarks. Overlapping 95% CIs of alt-nn1 and alt-nn3 ([49.4, 53.6] vs. [46.3, 59.9]) are due to severe sample size disparity (1,268 vs. 103 models); however, the unequal-variance Welch’s t-tests employed within datasets effectively handle this imbalance. Detailed per-dataset performance is available in **Section S4**.

Variant	n	Models	Mean (%)	SD (%)	95 % CI
alt-nn1	1	1,268	51.5	29.5	[49.4, 53.6]
alt-nn2	2	306	49.8	32.3	[45.1, 54.5]
alt-nn3	3	103	53.1	26.9	[46.3, 59.9]
alt-nn4	4	102	47.3	32.5	[39.0, 55.6]
alt-nn5	5	121	43.0	33.1	[35.3, 50.7]

Key Observations: The results reveal a clear three-zone regime. In the stable zone at $n=3$, the highest balanced mean of 53.1% is observed (+1.6% over $n=1$), with statistically significant improvement on fine-grained tasks (CIFAR-100: $p=0.001$, $d=0.73$). In the diminishing-returns zone at $n=4$ and $n=5$, accuracy drops to 47.3% and 43.0% respectively, as heterogeneous signals from additional examples begin to compete rather than complement. Finally, the collapse zone at $n=6$ produces only 7 valid models out of 3,394 total queries (99.8% failure rate), indicating the prompt has exceeded the model’s effective context capacity — the input examples consume the token budget, leaving insufficient space for a complete, syntactically valid output. This three-zone pattern is consistent with few-shot learning theory [2, 21], but the transition is sharper here because neural architecture code is substantially longer than natural language demonstrations, compressing the useful context window and making saturation both earlier and

414 more severe.

415 While absolute differences are modest under the rapid
416 screening protocol, the consistent trend across datasets in-
417 dicates that moderate context enrichment provides a prac-
418 tical balance between diversity and stability in LLM-based
419 architecture synthesis.

420 **The prompt structure follows this template:** Our few-
421 shot architecture prompting explicitly provides the refer-
422 ence architecture, the requested supporting model exam-
423 ples, performance-based accuracy guidance, and strict gen-
424 eration constraints. A complete, unredacted version of our
425 exact prompt template is provided in **Section S2 of the Sup-**
426 **plementary Material.**

427 **Prompt Design for Architecture Merging** The key
428 insight of our prompting strategy is the explicit instruction
429 to “combine best elements” and “combine best features,”
430 which directs the LLM toward architectural synthesis rather
431 than simple replication. Each supporting model is presented
432 with its achieved accuracy, providing performance context
433 that guides the merging process.

435 **Examples of Merged Architectural Concepts from** 436 **Generated Models**

437 Analysis of our generated `alt-nn3` architectures revealed
438 several cases where the LLM successfully synthesized
439 ideas from multiple computer vision models. Key ex-
440 amples include merging ResNet-style residual connections
441 with AlexNet-like classifiers, synthesizing Dual Path Net-
442 work (DPN) bottleneck blocks with progressive convolu-
443 tion backbones, and combining hierarchical residual units
444 with multi-scale feature pipelines. We provide full PyTorch
445 source code and detailed architectural analysis for these
446 synthesized models in **Section S3 of the Supplementary**
447 **Material.**

448 **Contrast with Single Supporting Example ($n = 1$)**
449 Analysis of representative `alt-nn1` variants reveals a con-
450 sistent pattern of shallow structural variation that con-
451 trasts sharply with the synthesis enabled by $n=3$. All ex-
452 amined $n=1$ models follow simple AlexNet-inspired se-
453 quential pipelines (`Conv`→`ReLU`→`Pool`) with only minor
454 channel-count mutations, suggesting the LLM defaults to
455 copying the reference architecture with superficial pertur-
456 bations rather than exploring the design space. No modular
457 abstractions appear (none of the `AirUnit` or `DPNBlock`-
458 style reusable components observed in `alt-nn3` are present)
459 and modern patterns such as residual connections, batch
460 normalisation in convolutional layers, and dual-path feature
461 fusion are entirely absent. This homogeneity confirms that
462 single-example prompting anchors the LLM too strongly to
463 the reference template, suppressing the generative diversity
464 that makes LLM-based NAS attractive, and directly moti-
465 vates the $n=3$ setting as the minimum context size needed
466 to break this anchoring effect.

5.2.1 **Code Illustration: Hybrid Architecture Synthe-** 467 **sis** 468

469 Figure 3 presents a qualitative comparison of archi-
470 tectures from our experiments, concretely demon-
471 strating the structural difference between single-
472 and three-example prompting. The `alt-nn1` variant
473 (`0e40be6fbc3426f57a305bfd8b8148fa`) fol-
474 lows an AlexNet-inspired sequential pipeline with
475 GELU activations and cascading max pooling, char-
476 acteristic of the shallow variation pattern observed
477 across $n=1$ outputs. In contrast, the `alt-nn3` variant
478 (`34df74344dd63a558c4b6413b809f6ed`) syn-
479 thesises ResNet-style residual units (`AirUnit` blocks
480 with identity shortcuts and BatchNorm) with a massive
481 fully-connected classifier head (4096 units with dropout)
482 — a cross-paradigm hybrid pattern entirely absent in $n=1$
483 variants and arising only when the LLM has multiple
484 architectural examples to draw from.

485 **Prompt Engineering Insight.** The explicit
486 `IMPROVEMENT RULES` constraints ensure syntactic
487 consistency while preserving the LLM’s freedom to inno-
488 vate in architectural design. Presenting accuracy metrics
489 alongside each supporting model provides performance-
490 based guidance, encouraging the LLM to favour patterns
491 from higher-performing architectures during synthesis
492 rather than averaging indiscriminately across examples.

5.3. **Comparison with Prior LLM-Based NAS** 493 494 **Methods**

495 **Methodological Context.** Direct quantitative comparison
496 with EvoPrompting [3] and LLMatic [22] is challenging
497 due to different evaluation protocols: EvoPrompting eval-
498 uates on MNIST-1D and CLRS (specialized benchmarks),
499 LLMatic on CIFAR-10 with NAS-Bench-201 (architectural
500 quality metrics), while we focus on standard computer vi-
501 sion benchmarks with rapid screening. However, we can
502 compare generation scale and qualitative insights.

503 **Generation Scale Comparison.**
504 • EvoPrompting: ~ 200 architectures across 2 benchmarks
505 • LLMatic: $\sim 2,000$ architectures on CIFAR-10 and NAS-
506 Bench-201
507 • Our work: 1,900 architectures across 7 diverse vision
508 benchmarks

509 **Complementary Insights.** While EvoPrompting
510 demonstrates that evolutionary prompt optimization can im-
511 prove generation quality, our finding that $n = 3$ provides
512 stable performance suggests that even simple random sam-
513 pling of examples achieves competitive results without iter-
514 ative refinement. This indicates that *strategic context en-*
515 *richment* may be as important as sophisticated selection
516 mechanisms.

517 LLMatic’s quality-diversity approach generates archi-
518 tectures with explicit Pareto trade-offs (accuracy vs. pa-

alt-nn1 (Single Example, $n=1$)

```
self.features = nn.Sequential(
    nn.Conv2d(in_shape[1], 64, kernel_size=11, stride=4,
              padding=2),
    nn.GELU(),
    nn.MaxPool2d(kernel_size=3, stride=2),
    nn.Conv2d(64, 192, kernel_size=5, padding=2),
    nn.GELU(),
    nn.MaxPool2d(kernel_size=3, stride=2))
```

alt-nn3 (Three Examples, $n=3$)

```
class AirUnit(nn.Module):
    def __init__(self, in_ch, out_ch, stride=1):
        # ... Init convolutional and downsample layers
        ...
    def forward(self, x):
        identity = x
        out = F.relu(self.bn1(self.conv1(x)))
        out = self.bn2(self.conv2(out))
        if self.down: identity = self.down(x)
        return F.relu(out + identity)
```

Figure 3. Overview of real generated architectures from our experiments. **Top:** alt-nn1 (0e40be6f...) — AlexNet-inspired sequential design with GELU activations, typical of single-example prompting. **Bottom:** alt-nn3 (34df7434...) — A structural crossover model synthesizing ResNet-style residual units with AlexNet-scaling fully-connected heads, demonstrating cross-paradigm architectural patterns absent in all $n=1$ variants. Provide full PyTorch source definitions in **Section S3 of the Supplementary Material**.

rameters). Our work is complementary: optimal few-shot prompting ($n = 3$) could enhance LLMatic’s generation phase by providing better architectural starting points for QD search.

Extension of NNGPT. Our work builds directly upon NNGPT [12], the base framework we use for generation. While NNGPT demonstrates powerful one-shot generation capabilities with single-reference prompting ($n = 1$), we demonstrate that enriching prompts with $n = 3$ supporting examples improves balanced accuracy by +1.6% on our rapid screening protocol. This suggests that prompt engineering can enhance NNGPT’s generation quality without architectural changes to the underlying system.

5.4. Per-Dataset and Statistical Analysis

While detailed per-dataset performance and statistical significance tests across all six computer vision benchmarks are provided in **Section S4**, key highlights demonstrate the task-dependency of context enrichment. Moderate context enrichment ($n=3$) is particularly beneficial in fine-grained classification settings (e.g., CIFAR-100: +11.6% vs $n=1$, $p = 0.001$, Cohen’s $d = 0.73$). CelebA-Gender also showed improvement (up to +6.5% for $n=2$, $p = 0.038$). In contrast, configurations with $n > 3$ show statistically significant lower performance on structured datasets like Ima-

geNette (−14.5%, $p = 0.010$) and CIFAR-10 (−8.4%, $p = 0.016$), indicating that excessively large prompt contexts can introduce instability in early training performance.

5.5. Computational Efficiency

Hash Validation Performance:

Our whitespace-normalized hashing approach provides substantial computational advantages over traditional code deduplication methods:

- **Computation time:** <1ms per code sample (measured on 4,033 generated architectures)
- **Baseline comparison:** AST parsing requires 10-100ms per sample
- **Speedup:** 100× faster than structural comparison
- **Accuracy:** Zero false positives across all generated code
- **Scalability:** Sub-millisecond latency enables real-time integration into generation loops

Practical Impact:

The computational efficiency of our approach has two critical implications for large-scale LLM-based architecture search:

(1) **Real-time duplicate detection:** Unlike AST-based methods that introduce 10-100ms overhead per architecture, our <1 ms validation enables duplicate checking within the generation loop without perceptible latency. This is essential for interactive exploration and iterative refinement workflows.

(2) **Prevention of redundant training:** In our experiments generating 1,900+ unique architectures, we observed that LLMs frequently produce formatting-level duplicates (identical code with different whitespace or indentation). Our normalization-based approach eliminates these redundant variants before training, preventing wasted computational resources. The elimination of redundant training runs directly reduces GPU utilization proportional to the duplicate rate observed during generation.

Comparison to Alternative Approaches:

Table 2 compares our method against existing code deduplication techniques.

Our approach optimally targets the dominant duplication pattern in LLM generation—formatting variations—while maintaining the computational efficiency of simple hashing. For applications requiring semantic equivalence detection, our method could serve as a fast first-pass filter followed by selective AST-based validation.

5.6. Discussion

Empirical Positioning. Our findings complement and extend existing LLM-based NAS work:

vs. EvoPrompting [3]: While evolutionary prompt optimization achieves 99.2% on MNIST-1D, our results on standard MNIST (> 96% across all variants) confirm that

Table 2. Code deduplication method comparison. ✓ = supported; ✗ = not supported.

Method	Time	Fmt.	Sem.
Raw string hash	<1ms	✗	✗
Winnowing [29]	2–5ms	Partial	✗
AST parsing [28]	10–100ms	✓	Partial
GraphCodeBERT [10]	50–200ms	✓	✓
Our approach	<1ms	✓	✗

strategic few-shot prompting is competitive without iterative refinement. The key difference: EvoPrompting requires costly prompt-tuning across generations; our $n = 3$ strategy works out-of-the-box.

vs. LLMatic [22]: Both approaches generate $\sim 2,000$ architectures, but target different objectives. LLMatic optimizes for multi-objective diversity (accuracy vs. size); we investigate prompt configuration effects. These are complementary: combining optimal few-shot prompting ($n = 3$) with QD search could yield both better starting architectures and more efficient exploration.

Extension of NNGPT [12]: We demonstrate that NNGPT’s generation quality can be enhanced through prompt engineering alone, without modifying the base framework. This “drop-in” improvement (+1.6% balanced accuracy, $p < 0.001$ on CIFAR-100) suggests that prompt optimization is an underexplored dimension of LLM-based architecture search.

Why $n=3$ is Optimal. Three supporting examples occupy the sweet spot between architectural diversity and prompt coherence. Too few examples ($n \leq 2$) produce outputs that replicate the reference model with minor mutations, as shown in Section 5.2.1. Beyond $n=3$, three compounding failure modes emerge: attention is spread too thinly across heterogeneous examples (context dilution); incompatible architectural choices introduce ambiguous generation signals (pattern conflict); and insufficient token budget remains for a complete output (context exhaustion). Each architectural example consumes ~ 500 tokens — an order of magnitude more than a natural language demonstration — making saturation both earlier and more catastrophic than in standard few-shot settings [2].

Task Complexity, Context Overflow, and Evaluation. The +11.6% gain on CIFAR-100 versus negligible effect on MNIST confirms that few-shot context benefits scale with task complexity: fine-grained 100-class classification requires the expressive feature hierarchies that only multi-example synthesis produces. The $n=6$ collapse (99.8% failure) is itself a valuable finding, establishing a concrete empirical upper bound on prompt context for this generation setting. Finally, our dataset-balanced macro-averaging prevents evaluation bias of up to 15.7% that arises from un-

weighted aggregation when generation success rates vary across datasets — a methodological concern applicable to any heterogeneous LLM-based NAS evaluation, and one we recommend as standard practice for future work in this area.

6. Conclusion

We present the first systematic empirical study of few-shot example count scaling in LLM-driven neural architecture generation. Building on NNGPT [12], we rigorously investigate how varying in-context example counts influences generation stability, architectural diversity, and early-epoch performance across six computer vision benchmarks.

Key Findings and the Stable Regime. Across 1,900 generated architectures, we identified that three supporting examples ($n = 3$) consistently provide a highly effective performance regime. This configuration achieves a +1.6% balanced mean accuracy improvement over standard single-example prompting, with statistically significant gains on complex tasks like CIFAR-100 (+11.6%, $p=0.001$). Qualitative assessments confirm $n=3$ uniquely enables sophisticated structural synthesis—such as merging residual pipelines with fully-connected layers—absent in $n=1$ variants. Conversely, extending beyond $n=3$ demonstrates severe diminishing returns, culminating in a 99.8% failure rate at $n=6$, establishing a definitive upper bound on prompt capacity.

Practical Contributions. To support large-scale generation, our highly optimized whitespace-normalized hashing technique yields a $100\times$ computational speedup over AST-based deduplication algorithms [28]. This enables real-time duplicate detection that prevents costly redundant training. Additionally, our dataset-balanced evaluation methodology effectively curtails statistical bias across highly heterogeneous benchmarks.

Positioning and Impact. Complementing recent advancements in evolutionary search [3] and quality-diversity algorithms [22], our work isolates the specific variable of prompt enrichment. We demonstrate that strategic prompt engineering alone ($n = 3$) substantially elevates the baseline quality of existing generative pipelines without complex, multi-stage selection mechanisms.

Limitations and Future Work. While single-epoch screening enables an unprecedented scale of comparative analysis, it inherently limits absolute precision. Future research must validate these findings using comprehensive multi-epoch training protocols and zero-cost proxies across broader vision tasks such as object detection and semantic segmentation. Developing hybrid pipelines that fuse our fast hashing with the deep semantic validation of AST-based models [10] also holds tremendous potential.

Ultimately, moderate few-shot enrichment makes LLM-based neural architecture search highly accessible and efficient for researchers on a budget.

687 **References**

- 688 [1] Mohamed S. Abdelfattah, Abhinav Mehrotra, Łukasz
689 Dudziak, and Nicholas D. Lane. Zero-cost proxies for
690 lightweight NAS. In *International Conference on Learning*
691 *Representations*, 2021. 4
- 692 [2] Tom B. Brown, Benjamin Mann, Nick Ryder, Melanie Sub-
693 biah, Jared Kaplan, Prafulla Dhariwal, et al. Language mod-
694 els are few-shot learners. *Advances in Neural Information*
695 *Processing Systems (NeurIPS)*, 33:1877–1901, 2020. 3, 5, 8
- 696 [3] Angelica Chen, David M. Dohan, and David R. So. Evo-
697 prompting: Language models for code-level neural architec-
698 ture search. In *NeurIPS*, 2023. 1, 2, 3, 6, 7, 8
- 699 [4] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan, Hen-
700 rique Ponde de Oliveira Pinto, Jared Kaplan, and Wojciech
701 Zaremba. Evaluating large language models trained on code.
702 *arXiv preprint arXiv:2107.03374*, 2021. 2
- 703 [5] Mark Chen, Jerry Tworek, Heewoo Jun, Qiming Yuan,
704 et al. Program synthesis with large language models. *arXiv*
705 *preprint arXiv:2108.07732*, 2022. 1
- 706 [6] Wuyang Chen, Xinyu Gong, and Zhangyang Wang. Neural
707 architecture search on ImageNet in four GPU hours: A theo-
708 retically inspired perspective. In *International Conference*
709 *on Learning Representations*, 2021. 4
- 710 [7] DeepSeek-AI. DeepSeek-Coder: When the large lan-
711 guage model meets programming. *arXiv preprint*
712 *arXiv:2401.14196*, 2024. 2, 3
- 713 [8] Romain Egele, Felix Mohr, Tom Viering, and Prasanna Bal-
714 aprakash. The unreasonable effectiveness of early discarding
715 after one epoch in neural network hyperparameter optimiza-
716 tion. *Neurocomputing*, 597:127964, 2024. 4
- 717 [9] Arash Torabi Goodarzi, Roman Kochnev, Waleed Khalid,
718 Furui Qin, Tolgay Atinc Uzun, Yashkumar Sanjaybhai
719 Dhameliya, Yash Kanubhai Kathiriyai, Zofia Antonina Ben-
720 tyn, Dmitry Ignatov, and Radu Timofte. Lemur neural net-
721 work dataset: Towards seamless automl. *arXiv preprint*
722 *arXiv:2504.10552*, 2025. 1, 2, 3
- 723 [10] Daya Guo, Shuo Ren, Shuai Lu, Zhangyin Feng, Duyu Tang,
724 Shujie Liu, Long Zhou, Nan Duan, Alexey Svyatkovskiy,
725 Shengyu Fu, et al. Graphcodebert: Pre-training code repre-
726 sentations with data flow. *arXiv preprint arXiv:2009.08366*,
727 2021. 3, 8
- 728 [11] Edward J. Hu, Yelong Shen, Phillip Wallis, Zeyuan Allen-
729 Zhu, Yuanzhi Li, Shean Wang, Lu Wang, and Weizhu Chen.
730 LoRA: Low-rank adaptation of large language models. In *Inter-
731 national Conference on Learning Representations (ICLR)*,
732 2022. 3
- 733 [12] Roman Kochnev, Waleed Khalid, Tolgay Atinc Uzun, Xi
734 Zhang, Yashkumar Sanjaybhai Dhameliya, Furui Qin, Chan-
735 dini Vysyaraju, Raghuvir Duvvuri, Avi Goyal, Dmitry Igna-
736 tov, and Radu Timofte. Nngpt: Rethinking automl with large
737 language models. *arXiv preprint arXiv:2511.20333*, 2025.
738 1, 2, 3, 7, 8
- 739 [13] Alex Krizhevsky. Learning multiple layers of features from
740 tiny images. Technical report, University of Toronto, 2009.
741 4
- 742 [14] Yann LeCun, Léon Bottou, Yoshua Bengio, and Patrick
743 Haffner. Gradient-based learning applied to document recog-
744 nition. *Proceedings of the IEEE*, 86(11):2278–2324, 1998.
745 4
- [15] Raymond Li, Loubna Ben Allal, Yangtian Zi, Niklas Muen-
746 nighoff, Denis Kocetkov, Chenghao Mou, and Leandro von
747 Werra. StarCoder: May the source be with you! *arXiv*
748 *preprint arXiv:2305.06161*, 2023. 2
- [16] Ming Lin, Pichao Wang, Zhenhong Sun, Heseng Chen, Xi-
749 yu Sun, Qi Qian, Hao Li, and Rong Jin. Zen-NAS: A
750 zero-shot NAS for high-performance image recognition. In
751 *Proceedings of the IEEE/CVF International Conference on*
752 *Computer Vision*, pages 12265–12275, 2021. 4
- [17] Hanxiao Liu, Karen Simonyan, and Yiming Yang. DARTS:
753 Differentiable architecture search. In *International Confer-
754 ence on Learning Representations (ICLR)*, 2019. 2
- [18] Ziwei Liu, Ping Luo, Xiaogang Wang, and Xiaoou Tang.
755 Deep learning face attributes in the wild. In *Proceedings*
756 *of the IEEE International Conference on Computer Vision*
757 *(ICCV)*, pages 3730–3738, 2015. 4
- [19] Yao Lu, Max Bartolo, Alastair Moore, Sebastian Riedel, and
758 Pontus Stenetorp. Fantastically ordered prompts and where
759 to find them: Overcoming few-shot prompt order sensitivity.
760 In *Proceedings of the 60th Annual Meeting of the Associa-
761 tion for Computational Linguistics (ACL)*, pages 8086–8098,
762 2022. 1, 3, 4
- [20] Joe Mellor, Jack Turner, Amos Storkey, and Elliot J. Crow-
763 ley. Neural architecture search without training. In *Proce-
764 edings of the 38th International Conference on Machine Learn-
765 ing*, pages 7588–7598. PMLR, 2021. 4
- [21] Sewon Min, Xinxi Lyu, Ari Holtzman, Mikel Artetxe, Mike
766 Lewis, Hannaneh Hajishirzi, and Luke Zettlemoyer. Re-
767 thinking the role of demonstrations: What makes in-context
768 learning work? *arXiv preprint arXiv:2202.12837*, 2022. 1,
769 3, 4, 5
- [22] Umair Muhammad Nasir, Sam Earle, Julian Togelius, Steven
770 James, and Christopher W. Cleghorn. Llmatic: Neural archi-
771 tecture search via large language models and quality diver-
772 sity optimization. *arXiv preprint arXiv:2306.01102*, 2024.
773 1, 2, 3, 6, 8
- [23] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bis-
774 saccio, Bo Wu, and Andrew Y. Ng. Reading digits in nat-
775 ural images with unsupervised feature learning. *NIPS Work-
776 shop on Deep Learning and Unsupervised Feature Learning*,
777 2011. 4
- [24] Erik Nijkamp, Bo Pang, Hiroaki Hayashi, Lifu Tu, Huan
778 Wang, Yingbo Zhou, and Caiming Xiong. CodeGen: An
779 open large language model for code with multi-turn program
780 synthesis. *arXiv preprint arXiv:2203.13474*, 2023. 2
- [25] Hieu Pham, Melody Guan, Barret Zoph, Quoc V. Le, and
781 Jeff Dean. Efficient neural architecture search via parameter
782 sharing. In *International Conference on Machine Learning*
783 *(ICML)*, pages 4095–4104, 2018. 2
- [26] Esteban Real, Alok Aggarwal, Yanping Huang, and Quoc V.
784 Le. Regularized evolution for image classifier architecture
785 search. In *Proceedings of the AAAI Conference on Artificial*
786 *Intelligence*, pages 4780–4789, 2019. 1, 2
- [27] Baptiste Rozière, Jonas Gehring, Fabian Gloeckle, Sten
787 Sootla, Itai Gat, Xiaoqing Tan, and Gabriel Synnaeve. Code
788 789 790 791 792 793 794 795 796 797 798 799 800

- 801 Llama: Open foundation models for code. *arXiv preprint*
802 *arXiv:2308.12950*, 2023. 2
- 803 [28] Hitesh Sajnani, Vaibhav Saini, Jeffrey Svajlenko, Chan-
804 chal K. Roy, and Cristina V. Lopes. Sourcerercc: Scaling
805 code clone detection to big-code. In *Proceedings of the 38th*
806 *International Conference on Software Engineering (ICSE)*,
807 pages 1157–1168, 2016. 3, 8
- 808 [29] Saul Schleimer, Daniel S. Wilkerson, and Alex Aiken. Win-
809 nowing: Local algorithms for document fingerprinting. In
810 *Proceedings of the ACM SIGMOD International Conference*
811 *on Management of Data*, pages 76–85. ACM, 2003. 3, 8
- 812 [30] Jason Wei, Xuezhi Wang, Dale Schuurmans, Maarten
813 Bosma, Brian Ichter, Fei Xia, Ed Chi, Quoc Le, and Denny
814 Zhou. Chain-of-thought prompting elicits reasoning in large
815 language models. *Advances in Neural Information Process-*
816 *ing Systems (NeurIPS)*, 35:24824–24837, 2022. 1, 3
- 817 [31] Colin White, Mahmoud Safari, Rhea Sukthanker, Binxin Ru,
818 Thomas Elsken, Arber Zela, Debadeepta Dey, and Frank
819 Hutter. Neural architecture search: Insights from 1000 pa-
820 pers. *arXiv preprint arXiv:2301.08727*, 2023. 4
- 821 [32] Bolei Zhou, Agata Lapedriza, Aditya Khosla, Aude Oliva,
822 and Antonio Torralba. Places: A 10 million image database
823 for scene recognition. In *IEEE Transactions on Pattern Anal-*
824 *ysis and Machine Intelligence (PAMI)*, pages 1452–1464,
825 2018. 4
- 826 [33] Barret Zoph and Quoc V. Le. Neural architecture search
827 with reinforcement learning. In *International Conference on*
828 *Learning Representations (ICLR)*, 2017. 1, 2

Enhancing LLM-Based Neural Network Generation: Few-Shot Prompting and Efficient Validation for Automated Architecture Design

Supplementary Material

829 S1. Detailed Algorithms

830 S1.1. Few-Shot Architecture Prompting (FSAP)

831 The FSAP procedure formalises the prompt construction
832 strategy described in Section 3.2 of the main paper. Given
833 a target dataset \mathcal{D} and desired context size n , the algorithm
834 queries LEMUR for high-performing architectures, desig-
835 nates one as the reference model, and draws n supporting
836 examples via uniform random sampling to isolate the effect
837 of context size from example-quality effects. The generated
838 prompt is passed to DeepSeek Coder 7B at temperature 0.6
839 to produce a complete PyTorch architecture.

Algorithm 1 Few-Shot Architecture Prompting (FSAP)

Require: Dataset \mathcal{D} , example count $n \in \{1, \dots, 6\}$

Require: LEMUR database \mathcal{L}

Ensure: Generated architecture code \mathcal{C}

```

1:  $\mathcal{M} \leftarrow \text{Query}(\mathcal{L}, \mathcal{D}, \text{top-accuracy})$ 
2:  $\mathcal{M}_{\text{ref}} \leftarrow \text{Sample}(\mathcal{M}, 1)$ 
3:  $\mathcal{M}_{\text{cand}} \leftarrow \mathcal{M} \setminus \mathcal{M}_{\text{ref}}$ 
4: if  $|\mathcal{M}_{\text{cand}}| \geq n$  then
5:    $\mathcal{M}_{\text{supp}} \leftarrow \text{RandomSample}(\mathcal{M}_{\text{cand}}, n)$ 
6: else
7:    $\mathcal{M}_{\text{supp}} \leftarrow \mathcal{M}_{\text{cand}}$ 
8: end if
9:  $\mathcal{P} \leftarrow \text{TaskDescription}(\mathcal{D})$ 
10:  $\mathcal{P} += \text{DatasetSpec}(\mathcal{D})$ 
11:  $\mathcal{P} += \text{FormatModel}(\mathcal{M}_{\text{ref}})$ 
12: for  $m \in \mathcal{M}_{\text{supp}}$  do
13:    $\mathcal{P} += \text{FormatSupporting}(m)$ 
14: end for
15:  $\mathcal{P} += \text{GenerationRules}()$ 
16:  $\mathcal{C} \leftarrow \text{DeepSeekCoder}(\mathcal{P}, T=0.6)$ 
17: return  $\mathcal{C}$ 

```

840 S1.2. Whitespace-Normalized Hash Validation

841 The hash validation procedure safely eliminates formatting-
842 level duplicate architectures prior to the computationally
843 demanding training phase, effectively saving 2–3 GPU
844 hours per rejected instance. The automated three-step pro-
845 cess — comprehensive whitespace removal, MD5 hashing,
846 and B-tree indexed database lookup — runs in $O(|\mathcal{C}| +$
847 $\log N)$ time. This is heavily dominated by the linear code-
848 length term, consistently achieving sub-millisecond latency
849 in practical applications.

Algorithm 2 Whitespace-Normalized Hash Validation

Require: Code string \mathcal{C} , LEMUR database \mathcal{L}

Ensure: Decision: {ACCEPT, REJECT}

```

1:  $\mathcal{C}' \leftarrow \text{RemoveWhitespace}(\mathcal{C})$   $\triangleright O(|\mathcal{C}|)$ 
2:  $h \leftarrow \text{MD5}(\mathcal{C}')$   $\triangleright O(|\mathcal{C}|), \text{hardware-accel.}$ 
3:  $\mathcal{H} \leftarrow \text{Query}(\mathcal{L}, \text{SELECT nn\_id FROM lemur})$   $\triangleright$   

 $O(\log N)$ 
4: if  $h \in \mathcal{H}$  then
5:   return REJECT
6: else
7:   return ACCEPT
8: end if

```

The total asymptotic complexity of $O(|\mathcal{C}| + \log N) \approx$ 850
 $O(|\mathcal{C}|)$ makes this hash validation significantly faster than 851
standard AST parsing ($O(|\mathcal{C}|^{1.5})$ in practice) or Graph- 852
CodeBERT inference (50–200ms per sample). Further- 853
more, this lightweight approach directly targets the most 854
dominant duplication pattern observed in LLM pipelines: 855
superficial formatting variations arising from inconsistent 856
indentation and redundant whitespace within raw code gen- 857
eration outputs. 858

S2. Complete Prompt Template 859

The full prompt template utilized across all FSAP experi- 860
ments is detailed below. The construction of this prompt 861
acts as the vital control interface for the LLM-driven gener- 862
ation pipeline. The key design choices driving this template 863
are: 864

(1) **Accuracy labels alongside each model:** By pro- 865
viding explicit performance metrics mapped to each code 866
block, we establish a quantifiable, performance-based guid- 867
ance signal. This context encourages the LLM to correlate 868
specific architectural motifs with high empirical success, bi- 869
asing its output toward proven structural patterns rather than 870
arbitrary mutations. 871

(2) **Explicit synthesis instructions:** The directive to 872
“combine best elements” deliberately shifts the model’s 873
generative priority away from simply duplicating the refer- 874
ence architecture or making shallow alterations. Instead, 875
it actively promotes structural crossover, challenging the 876
LLM to intelligently synthesize novel, hybrid architectures 877
by merging advantageous features from the diverse support- 878
ing examples. 879

(3) **Strict output constraints:** By strictly enforcing 880
a fixed class name (`Net`), rigid method signatures, and 881
PyTorch-only dependencies, we provide a reliable scaffold 882

883 that guarantees every generated architecture remains syn-
 884 tactically valid. This ensures zero friction when integrating
 885 models into our automated training pipeline, preserving the
 886 LLM’s architectural freedom in designing internal topolo-
 887 gies.

```

888 CREATE an IMPROVED neural network by combining the best
889 elements from these models:
890
891 MAIN MODEL (current accuracy: {accuracy}%):
892 ```python
893 {reference_architecture_code}
894 ```
895
896 SUPPORTING MODEL 1 (accuracy: {addon_accuracy_1}%):
897 ```python
898 {supporting_architecture_1_code}
899 ```
900
901 SUPPORTING MODEL 2 (accuracy: {addon_accuracy_2}%):
902 ```python
903 {supporting_architecture_2_code}
904 ```
905
906 [... up to n supporting models ...]
907
908 IMPROVEMENT RULES - FOLLOW EXACTLY:
909 1. Class name: 'Net' (unchanged)
910 2. Methods: __init__, forward, train_setup(device),
911    learn(data,target,device) - keep signatures
912 3. Include: supported_hyperparameters() - ['lr', '
913    momentum']
914 4. Only standard PyTorch (no torchvision)
915 5. Works with 32x32 RGB images - [num_classes] classes
916
917 IMPROVE by combining best features from all models above
918
919 Provide COMPLETE improved model code:
920
  
```

922 S3. Complete Architecture Code Examples

923 This section provides the full PyTorch source code for the
 924 four representative architectures discussed qualitatively in
 925 Section 5.2.1 of the main paper. The first three (alt-nn3 vari-
 926 ants) illustrate the architectural synthesis enabled by $n=3$
 927 prompting; the fourth (alt-nn1 baseline) demonstrates the
 928 shallow variation characteristic of single-example prompt-
 929 ing for direct comparison.

930 S3.1. Example 1: ResNet + AlexNet Hybrid

931 **Model ID:** *alt-nn3-34df74344dd63a558c4b6413b809f6ed*

932 This architecture exemplifies successful cross-paradigm
 933 synthesis, seamlessly integrating modern representation
 934 learning techniques with classical, high-capacity classifica-
 935 tion heads. Specifically, it synthesises ResNet-style resid-
 936 ual units (AirUnit blocks featuring identity skip connec-
 937 tions and batch normalisation to mitigate vanishing gradi-
 938 ents) with AlexNet’s signature large-scale, two-stage fully-
 939 connected classifier (utilizing 4096-unit dense layers regu-
 940 larized by aggressive dropout).

941 This structural hybrid leverages the deep, stable feature
 942 extraction capabilities of residual networks while preserv-
 943 ing the massive parameterization of early deep learning

classifiers. Crucially, such sophisticated topological com-
 binations are qualitatively absent in all $n=1$ generated vari-
 ants, which typically default to shallow, sequential template
 alterations. This strongly corroborates that providing three
 supporting examples uniquely enables the LLM to perform
 genuine architectural pattern merging across disparate net-
 work families.

```

import torch
import torch.nn as nn
import torch.nn.functional as F

class AirUnit(nn.Module):
    """ResNet-style residual block"""
    def __init__(self, in_ch, out_ch, stride=1):
        super().__init__()
        self.conv1 = nn.Conv2d(in_ch, out_ch, 3, stride,
            1, bias=False)
        self.bn1 = nn.BatchNorm2d(out_ch)
        self.conv2 = nn.Conv2d(out_ch, out_ch, 3, 1, 1,
            bias=False)
        self.bn2 = nn.BatchNorm2d(out_ch)
        self.down = None
        if stride != 1 or in_ch != out_ch:
            self.down = nn.Sequential(
                nn.Conv2d(in_ch, out_ch, 1, stride, bias
                    =False),
                nn.BatchNorm2d(out_ch))

    def forward(self, x):
        identity = x
        out = F.relu(self.bn1(self.conv1(x)))
        out = self.bn2(self.conv2(out))
        if self.down:
            identity = self.down(x)
        return F.relu(out + identity)

class Net(nn.Module):
    def __init__(self, num_classes=10):
        super().__init__()
        self.l1 = AirUnit(3, 64, stride=2)
        self.l2 = AirUnit(64, 128, stride=2)
        self.l3 = AirUnit(128, 256, stride=2)
        self.pool = nn.AdaptiveAvgPool2d((1,1))
        self.fc1 = nn.Linear(256, 4096) # AlexNet
        self.drop = nn.Dropout(0.5)
        self.fc2 = nn.Linear(4096, 4096)
        self.out = nn.Linear(4096, num_classes)

    def forward(self, x):
        x = self.l3(self.l2(self.l1(x)))
        x = torch.flatten(self.pool(x), 1)
        x = self.drop(F.relu(self.fc1(x)))
        x = self.drop(F.relu(self.fc2(x)))
        return self.out(x)

    def supported_hyperparameters(self):
        return ['lr', 'momentum']
  
```

S3.2. Example 2: DPN-Inspired Hybrid

Model ID: *alt-nn3-57d770565afc5e0d651cb0938fc8f942*

This specific generated architecture successfully merges
 standard Dual Path Network (DPN) bottleneck feature-
 extraction blocks($1\times 1 \rightarrow 3\times 3 \rightarrow 1\times 1$) with a highly
 progressive and customized convolutional backbone design.
 The unusually asymmetric mid-network channel progres-
 sion ($64\rightarrow 128\rightarrow 440\rightarrow 384\rightarrow 192$) serves as strong empir-
 ical evidence of the model’s creative architectural synthesis
 rather than mere direct template copying. Crucially, this

exact dimensional configuration appears in absolutely none of the provided LEMUR reference models utilized in the initial prompt context. Furthermore, the final classification head effectively incorporates batch normalization directly within the fully-connected dense layers—a distinctly modern architectural regularization pattern that clearly demonstrates robust, cross-architecture feature pattern transfer during the generation phase.

```

import torch
import torch.nn as nn
import torch.nn.functional as F

class DPNBlock(nn.Module):
    """DPN-inspired bottleneck"""
    def __init__(self, in_ch, out_ch, stride=1):
        super().__init__()
        self.c1 = nn.Conv2d(in_ch, out_ch, 1, bias=False)
        self.b1 = nn.BatchNorm2d(out_ch)
        self.c2 = nn.Conv2d(out_ch, out_ch, 3, stride,
1, bias=False)
        self.b2 = nn.BatchNorm2d(out_ch)
        self.c3 = nn.Conv2d(out_ch, out_ch, 1, bias=
False)
        self.b3 = nn.BatchNorm2d(out_ch)
        self.sc = nn.Sequential()
        if stride != 1 or in_ch != out_ch:
            self.sc = nn.Sequential(
                nn.Conv2d(in_ch, out_ch, 1, stride, bias
=False),
                nn.BatchNorm2d(out_ch))

    def forward(self, x):
        out = F.relu(self.b1(self.c1(x)))
        out = F.relu(self.b2(self.c2(out)))
        out = self.b3(self.c3(out))
        return F.relu(out + self.sc(x))

class Net(nn.Module):
    def __init__(self, num_classes=10):
        super().__init__()
        # Unusual 440-channel mid-layer: creative
        synthesis
        self.l1 = DPNBlock(3, 64, 1)
        self.l2 = DPNBlock(64, 128, 2)
        self.l3 = DPNBlock(128, 440, 2)
        self.l4 = DPNBlock(440, 384, 1)
        self.l5 = DPNBlock(384, 192, 2)
        self.pool = nn.AdaptiveAvgPool2d((1,1))
        self.fc1 = nn.Linear(192, 512)
        self.bn = nn.BatchNorm1d(512)
        self.fc2 = nn.Linear(512, num_classes)

    def forward(self, x):
        x = self.l5(self.l4(self.l3(self.l2(self.l1(x))))
        x = torch.flatten(self.pool(x), 1)
        return self.fc2(F.relu(self.bn(self.fc1(x))))

    def supported_hyperparameters(self):
        return ['lr', 'momentum']

```

S3.3. Example 3: Hierarchical Residual Units + Multi-Scale Features

Model ID: *alt-nn3-dcb59f747eb3b27c0552d2aea356e33a*

This architecture synthesizes ResNet-style AirUnit blocks with three-layer residual paths into a hierarchical feature extraction pipeline combining varying spatial scales. The model begins with large kernel convolutions (7×7 ,

stride 3) for rapid downsampling before cascading into sophisticated residual units integrating pooling steps implicitly through strided convolutions. This design merges classical and modern approaches: aggressive early spatial reduction from AlexNet with deep residual learning from ResNet, creating a compact yet expressive architecture suitable for CIFAR-scale tasks.

```

class AirInitBlock(nn.Module):
    def __init__(self, in_channels, out_channels):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Conv2d(in_channels, out_channels,
kernel_size=3, stride=2, padding=1),
            nn.BatchNorm2d(out_channels),
            nn.ReLU(inplace=True)
        )
    def forward(self, x): return self.layers(x)

class AirUnit(nn.Module):
    def __init__(self, in_channels, out_channels, stride):
        super().__init__()
        self.layers = nn.Sequential(
            nn.Conv2d(in_channels, out_channels,
kernel_size=3, stride=stride, padding=1),
            nn.BatchNorm2d(out_channels), nn.ReLU(
inplace=True),
            nn.Conv2d(out_channels, out_channels,
kernel_size=3, stride=1, padding=1),
            nn.BatchNorm2d(out_channels), nn.ReLU(
inplace=True),
            nn.Conv2d(out_channels, out_channels,
kernel_size=3, stride=1, padding=1),
        )
        self.downsample = (
            nn.Sequential(
                nn.Conv2d(in_channels, out_channels,
kernel_size=1, stride=stride, bias=False),
                nn.BatchNorm2d(out_channels)
            ) if stride != 1 or in_channels !=
out_channels else nn.Identity()
        )
        self.relu = nn.ReLU(inplace=True)

    def forward(self, x):
        return self.relu(self.layers(x) + self.
downsample(x))

class Net(nn.Module):
    def __init__(self, in_shape: tuple, out_shape: tuple
, prm: dict, device: torch.device) -> None:
        super().__init__()
        self.features = nn.Sequential(
            nn.Conv2d(in_shape[1], 96, kernel_size=7,
stride=3, padding=2),
            nn.BatchNorm2d(96), nn.ReLU(inplace=True),
            nn.MaxPool2d(kernel_size=2, stride=2),
            AirInitBlock(96, 192),
            AirUnit(192, 384, stride=2),
            AirUnit(384, 256, stride=1),
            AirUnit(256, 256, stride=2),
            nn.AdaptiveAvgPool2d((6, 6))
        )
        self.classifier = nn.Sequential(
            nn.Dropout(p=prm['dropout']),
            nn.Linear(256 * 6 * 6, 4096), nn.ReLU(
inplace=True),
            nn.Dropout(p=prm['dropout']),
            nn.Linear(4096, out_shape[0])
        )

    def forward(self, x: torch.Tensor) -> torch.Tensor:
        return self.classifier(torch.flatten(self.
features(x), 1))

```

1159 **S3.4. Example 4: Baseline alt-nn1 (Single Example)**1160 **Model ID:** *alt-nn1-0e40be6fbc3426f57a305bfd8b8148fa*1161 This representative $n=1$ baseline architecture illustrates
1162 the shallow sequential variation pattern typical in single-
1163 prompt experiments. Lacking modular abstraction, residual
1164 shortcuts, or normalization, the model relies on a dense lin-
1165 ear progression of GELU-activated pooling operations char-
1166 acteristic of straightforward, non-creative structural muta-
1167 tion relative to references.

```

1168 class Net(nn.Module):
1169     def __init__(self, in_shape: tuple, out_shape: tuple
1170     , prm: dict, device: torch.device) -> None:
1171         super().__init__()
1172         self.features = nn.Sequential(
1173             nn.Conv2d(in_shape[1], 64, kernel_size=11,
1174             stride=4, padding=2),
1175             nn.GELU(),
1176             nn.MaxPool2d(kernel_size=3, stride=2),
1177             nn.Conv2d(64, 192, kernel_size=5, padding=2)
1178         ,
1179             nn.GELU(),
1180             nn.MaxPool2d(kernel_size=3, stride=2),
1181             nn.Conv2d(192, 384, kernel_size=3, padding
1182             =1),
1183             nn.GELU(),
1184             nn.Conv2d(384, 256, kernel_size=3, padding
1185             =1),
1186             nn.GELU(),
1187             nn.Conv2d(256, 256, kernel_size=3, padding
1188             =1),
1189             nn.GELU(),
1190             nn.MaxPool2d(kernel_size=3, stride=2)
1191         )
1192         self.avgpool = nn.AdaptiveAvgPool2d((6, 6))
1193         self.classifier = nn.Sequential(
1194             nn.Dropout(p=prm['dropout']),
1195             nn.Linear(256 * 6 * 6, 4096),
1196             nn.GELU(),
1197             nn.Dropout(p=prm['dropout']),
1198             nn.Linear(4096, 643),
1199             nn.GELU(),
1200             nn.Linear(643, out_shape[0])
1201         )
1202     def forward(self, x: torch.Tensor) -> torch.Tensor:
1203         return self.classifier(torch.flatten(self.
1204         avgpool(self.features(x)), 1))

```

1208 **S4. Extended Per-Dataset Results**1209 This section expands upon the dataset-balanced evaluation
1210 discussed in the main paper by providing the raw, per-
1211 dataset quantitative performance and statistical significance
1212 matrices.1213 **S4.1. Per-Dataset Performance Analysis**1214 Table S1 presents detailed per-dataset results with statistical
1215 significance markers.1216 **Observations:**

- 1217 •
- Task-dependent variation:**
- The relative performance of
-
- 1218 prompt variants differs across datasets. In particular,
- $n=3$
-
- 1219 achieves the highest mean accuracy on CIFAR-100, while
-
- 1220 other datasets exhibit different preferences.

Table S1. Per-Dataset Performance (Mean Accuracy %) on 1-
epoch. Asterisks denote significance vs. baseline: * $p < 0.05$,
** $p < 0.01$. Best per dataset in **bold**.

Dataset	alt-nn1	alt-nn2	alt-nn3	alt-nn4	alt-nn5
MNIST	96.5	93.8*	97.1	93.9	94.7
CelebA-Gender	75.8	82.3*	74.4	80.6	72.1
CIFAR-10	38.7	36.1	38.3	30.3*	34.0
CIFAR-100	14.5	7.4*	26.1**	10.8	10.5
ImageNette	44.2	36.4*	42.5	29.8**	18.2**
SVHN	39.2	43.1	40.0	38.3	28.5
Bal. Mean	51.5	49.8	53.1	47.3	43.0

- 1221 • **Simple vs. complex tasks:** On simpler benchmarks such
1222 as MNIST, performance remains high across all variants
1223 ($> 93\%$), suggesting limited sensitivity to prompt config-
1224 uration under low task complexity.
- 1225 • **Context scaling effects:** Variants with $n>3$ exhibit lower
1226 performance on several datasets (e.g., ImageNette), in-
1227 dicated that larger prompt contexts do not consistently
1228 translate into improved early performance.

1229 **S4.2. Statistical Significance Analysis**1230 Table S2 presents statistical validation focusing on key find-
1231 ings.Table S2. Statistical Significance Tests (per-dataset comparison
vs. alt-nn1 baseline). Only results with $p < 0.05$ shown.

Dataset	Comparison	Δ	p	d
CIFAR-100	alt-nn3 vs alt-nn1	+11.6%	0.001	0.73
MNIST	alt-nn2 vs alt-nn1	-2.7%	0.029	-0.19
CelebA	alt-nn2 vs alt-nn1	+6.5%	0.038	0.41
CIFAR-10	alt-nn4 vs alt-nn1	-8.4%	0.016	-0.54
ImageNette	alt-nn4 vs alt-nn1	-14.5%	0.010	-1.20

1232 **Key Statistical Result:** On CIFAR-100, the $n=3$
1233 configuration achieves higher mean accuracy than $n=1$
1234 ($p=0.001$, Cohen's $d=0.73$), suggesting that moderate con-
1235 text enrichment may be particularly beneficial in fine-
1236 grained classification settings under the rapid screening pro-
1237 tocol.1238 In contrast, configurations with $n>3$ show statistically
1239 significant lower performance on certain datasets (e.g., Im-
1240 ageNette), indicating that larger prompt contexts can in-
1241 troduce instability or diminishing returns in early training
1242 performance. We emphasise that these comparisons reflect
1243 early-epoch behavior under a rapid screening protocol and
1244 are intended to capture relative trends across prompt vari-
1245 ants rather than final converged performance.