

¹ LabelFusion: Learning to Fuse LLMs and Transformer Classifiers for Robust Text Classification

³ **Christoph Weisser^{1,2} and LabelFusion contributors¹**

⁴ 1 Campus-Institut Data Science, Göttingen, Germany 2 Centre for Statistics, Georg-August-Universität
⁵ Göttingen, Germany

DOI: [10.xxxxxx/draft](https://doi.org/10.xxxxxx/draft)

Software

- [Review ↗](#)
- [Repository ↗](#)
- [Archive ↗](#)

Editor: [Open Journals ↗](#)

Reviewers:

- [@openjournals](#)

Submitted: 01 January 1970

Published: unpublished

License

Authors of papers retain copyright and release the work under a Creative Commons Attribution 4.0 International License ([CC BY 4.0](#)).

⁶ Summary

⁷ LabelFusion is a fusion ensemble for text classification that learns to combine a traditional
⁸ transformer-based classifier (e.g., RoBERTa) with one or more Large Language Models (LLMs)
⁹ such as OpenAI GPT, Google Gemini, or DeepSeek to deliver accurate and cost-aware predictions
¹⁰ across multi-class and multi-label tasks. The package provides a simple high-level interface
¹¹ (`AutoFusionClassifier`) that trains the full pipeline end-to-end with minimal configuration,
¹² and a flexible API for advanced users. Under the hood, LabelFusion takes vector signals from
¹³ an ML backbone (logits) and LLM(s) (per-class scores), calibrates them using temperature
¹⁴ scaling and isotonic regression, and feeds their concatenation into a small multi-layer perceptron
¹⁵ (`FusionMLP`) that is trained to produce the final prediction. This learned fusion approach
¹⁶ captures complementary strengths of LLM reasoning and transformer efficiency, yielding robust
¹⁷ performance across domains—achieving 92.4% accuracy on AG News topic classification—while
¹⁸ enabling practical trade-offs between accuracy, latency, and cost.

Statement of Need

²⁰ Modern text classification spans diverse scenarios—from sentiment and topic tagging to policy
²¹ enforcement and routing—often under constraints that vary per deployment (throughput, cost
²² ceilings, data privacy). While transformer classifiers such as BERT/RoBERTa achieve strong
²³ supervised performance ([Devlin et al., 2018](#); [Liu et al., 2019](#)), frontier LLMs can excel in
²⁴ low-data, ambiguous, or cross-domain settings ([OpenAI, 2023](#)). No single model family is
²⁵ uniformly best: LLMs are powerful yet comparatively costly and rate-limited, whereas fine-tuned
²⁶ transformers are efficient but may struggle with out-of-distribution cases.

²⁷ LabelFusion addresses this gap by: (1) exposing a minimal “AutoFusion” interface that trains a
²⁸ learned combination of an ML backbone and one or more LLMs; (2) supporting both multi-class
²⁹ and multi-label classification; (3) providing calibration of LLM scores and ML logits for better
³⁰ probability estimates; and (4) integrating cleanly with existing ensemble utilities. Researchers
³¹ and practitioners can therefore leverage LLMs where they add value while retaining the speed
³² and determinism of transformer models.

³³ State of the Field

³⁴ Ensembles improve robustness by aggregating diverse predictors ([Dietterich, 2000](#); [Hansen &](#)
³⁵ [Salamon, 1990](#)). Mixture-of-experts approaches further specialize components and learn to
³⁶ combine their outputs ([Jacobs et al., 1991](#)). In applied NLP, common tools such as scikit-learn
³⁷ ([Pedregosa et al., 2011](#)) and Hugging Face Transformers ([Wolf et al., 2019](#)) offer strong
³⁸ baselines but do not provide a turnkey, learned fusion of LLMs with supervised transformers.
³⁹ Orchestration frameworks (e.g., LangChain) focus on tool use rather than classification
⁴⁰ ensembles. LabelFusion contributes a focused, production-minded implementation of a small
⁴¹ learned combiner that operates on calibrated per-class signals from both model families.

42 Functionality and Design

43 LabelFusion consists of three layers:

- 44 ▪ ML component: a RoBERTa-style classifier produces per-class logits for input texts.
- 45 ▪ LLM component(s): provider-specific classifiers (OpenAI, Claude, Gemini, DeepSeek)
- 46 return per-class scores via prompting. Scores are cached to minimize API calls and cost.
- 47 ▪ Fusion component: a compact MLP concatenates ML logits and LLM scores and outputs
- 48 fused logits. The ML backbone is trained/fine-tuned with a small learning rate; the fusion
- 49 MLP uses a higher rate, enabling rapid adaptation without destabilizing the encoder.

50 Key features:

- 51 ▪ **Multi-class and multi-label support** with consistent data structures and unified training
- 52 pipeline.
- 53 ▪ **Calibration of model signals** using temperature scaling and isotonic regression ([Guo et](#)
- 54 [al., 2017; Zadrozny & Elkan, 2002](#)) for better probability estimates. LLM scores are
- 55 calibrated on validation data before fusion training.
- 56 ▪ **LLM response caching** with disk-based persistence reduces API costs by storing and
- 57 reusing predictions. Cache invalidation handles configuration changes automatically.
- 58 ▪ **Batched scoring** processes multiple texts efficiently with configurable batch sizes for both
- 59 ML tokenization and LLM API calls.
- 60 ▪ **Results management** via ResultsManager tracks experiments, stores predictions, com-
- 61 putes metrics, and enables reproducible research workflows.
- 62 ▪ **Flexible interfaces**: Command-line training via `train_fusion.py` with YAML configs for
- 63 research; or minimal AutoFusion API for quick deployment.
- 64 ▪ **Composable design**: LabelFusion can serve as a strong base learner in higher-level
- 65 ensembles (e.g., voting/weighted combinations of multiple fusion models).

66 Minimal Example (AutoFusion)

```
from textclassify import AutoFusionClassifier

config = {
    'llm_provider': 'deepseek',
    'label_columns': ['positive', 'negative', 'neutral']
}

clf = AutoFusionClassifier(config)
clf.fit(train_dataframe)           # trains ML backbone, gathers LLM scores, fits fusi
pred = clf.predict(["This is amazing!"]) # fused prediction
```

67 CLI and Configuration

68 Users can generate a starter config and train via the command line:

- 69 ▪ Create config: `python train_fusion.py --create-config fusion_config.yaml`
- 70 ▪ Train: `python train_fusion.py --config fusion_config.yaml`
- 71 ▪ Optional test data and output artifacts are also supported.

72 Quality Control

73 The repository includes unit and integration tests (see `tests/`) that validate configuration

74 handling, core types, and package integration. Fusion-specific logic is exercised in examples

75 and the CLI, which run end-to-end training with deterministic seeds where applicable.

76 Evaluation scripts (`tests/evaluation/`) provide comprehensive benchmarking on standard

77 datasets: - **AG News** ([Zhang et al., 2015](#)): 4-class topic classification with experiments

⁷⁸ across varying training data sizes (20%–100%) - **GoEmotions** ([Demszky et al., 2020](#)): 28-class
⁷⁹ multi-label emotion classification for validating multi-label fusion performance

⁸⁰ LLM scoring paths implement retries and disk caching; transformer training supports standard
⁸¹ sanity checks (overfit a small batch, reduced batch sizes for constrained hardware). Metrics
⁸² (accuracy/F1, per-label scores) are computed automatically and stored with run artifacts to
⁸³ facilitate regression tracking and reproducibility.

⁸⁴ Availability and Installation

⁸⁵ LabelFusion is distributed as part of the `textclassify` package under the MIT license and
⁸⁶ is available at <https://github.com/DataandAIResearch/LabelFusion>. The fusion components
⁸⁷ require Python 3.8+ and common scientific Python dependencies (PyTorch, transformers,
⁸⁸ scikit-learn, numpy, pandas, PyYAML). Optional plotting depends on matplotlib/seaborn.
⁸⁹ Installation and quick-start snippets are provided in the README and FUSION_README.md.

⁹⁰ Production-Ready Features

⁹¹ Beyond the core fusion methodology, LabelFusion includes features for practical deployment:

- ⁹² ▪ **LLM Response Caching:** Automatic caching of LLM predictions to disk reduces API
⁹³ costs and enables reproducible experiments. The cache system handles invalidation and
⁹⁴ supports multiple cache backends.
- ⁹⁵ ▪ **Results Management:** Built-in `ResultsManager` tracks experiments, stores predictions,
⁹⁶ and computes metrics automatically. Supports comparison across runs and configuration
⁹⁷ tracking.
- ⁹⁸ ▪ **Batch Processing:** Efficient batched scoring of texts with configurable batch sizes for
⁹⁹ both ML and LLM components.
- ¹⁰⁰ ▪ **Cost Monitoring:** Tracks API usage and estimated costs across LLM providers with
¹⁰¹ configurable budget limits.

¹⁰² Impact and Use Cases

¹⁰³ Empirical Performance

¹⁰⁴ LabelFusion has been evaluated on standard benchmark datasets to validate its effectiveness.
¹⁰⁵ Key findings demonstrate consistent improvements over individual model components:

¹⁰⁶ AG News Topic Classification

¹⁰⁷ Evaluation on the AG News dataset ([Zhang et al., 2015](#)) (4-class topic classification) with
¹⁰⁸ 5,000 test samples shows:

Training Data	Model	Accuracy	F1-Score	Precision	Recall
20% (800)	Fusion	92.2%	0.922	0.923	0.922
20% (800)	RoBERTa	89.8%	0.899	0.902	0.898
20% (800)	OpenAI	84.4%	0.844	0.857	0.844
40% (1,600)	Fusion	92.2%	0.922	0.924	0.922
40% (1,600)	RoBERTa	91.0%	0.911	0.913	0.910
40% (1,600)	OpenAI	84.4%	0.844	0.857	0.844
100% (4,000)	Fusion	92.4%	0.924	0.926	0.924
100% (4,000)	RoBERTa	92.2%	0.922	0.923	0.922
100% (4,000)	OpenAI	84.4%	0.844	0.857	0.844

¹⁰⁹ **Key Observations:** - Fusion consistently outperforms individual models across all training
¹¹⁰ data sizes - With only 20% training data, Fusion achieves 92.2% accuracy—matching its

111 performance with full data - Demonstrates superior **data efficiency**: fusion learning extracts
112 maximum value from limited examples - RoBERTa alone requires 100% of data to approach
113 Fusion's 20% performance - LLM (OpenAI) shows stable but lower performance, highlighting
114 the value of combining approaches

115 These results validate that learned fusion captures complementary strengths: the LLM provides
116 robust reasoning even with limited training data, while the ML backbone adds efficiency and
117 domain-specific patterns.

118 Application Domains

119 Learned fusion excels in scenarios where model strengths complement each other:

- 120 ■ **Customer feedback analysis** with nuanced multi-label taxonomies where LLMs handle
121 ambiguous sentiment while ML models efficiently process clear cases
- 122 ■ **Content moderation** where uncertain cases benefit from LLM reasoning while routine
123 items rely on the fast ML backbone, enabling real-time processing with accuracy
124 guarantees
- 125 ■ **Scientific literature classification** across heterogeneous topics where domain shift is
126 common and LLMs provide robustness to new terminology
- 127 ■ **Low-resource settings** where limited training data is available but task complexity requires
128 sophisticated reasoning

129 The approach enables pragmatic cost control (e.g., the fusion layer learns when to rely more
130 heavily on the efficient ML backbone versus the more expensive LLM signal) while retaining a
131 single trainable decision surface that optimizes for the specific deployment constraints.

132 Acknowledgements

133 We thank contributors and users who reported issues and shared datasets. LabelFusion
134 builds on the open-source ecosystem, notably Hugging Face Transformers ([Wolf et al., 2019](#)),
135 scikit-learn ([Pedregosa et al., 2011](#)), PyTorch ([Paszke et al., 2019](#)), and LLM provider SDKs.
136 We acknowledge the use of the AG News and GoEmotions benchmark datasets for evaluation.

137 References

- 138 Demszky, D., Movshovitz-Attias, D., Ko, J., Cowen, A., Nemade, G., & Ravi, S. (2020).
139 GoEmotions: A dataset of fine-grained emotions. *Proceedings of the 58th Annual Meeting
140 of the Association for Computational Linguistics*, 4040–4054.
- 141 Devlin, J., Chang, M.-W., Lee, K., & Toutanova, K. (2018). BERT: Pre-training of deep
142 bidirectional transformers for language understanding. *arXiv Preprint arXiv:1810.04805*.
- 143 Dietterich, T. G. (2000). Ensemble methods in machine learning. *International Workshop on
144 Multiple Classifier Systems*, 1–15.
- 145 Guo, C., Pleiss, G., Sun, Y., & Weinberger, K. Q. (2017). On calibration of modern neural
146 networks. *Proceedings of the 34th International Conference on Machine Learning*.
- 147 Hansen, L. K., & Salamon, P. (1990). Neural network ensembles. *IEEE Transactions on
148 Pattern Analysis and Machine Intelligence*, 12(10), 993–1001.
- 149 Jacobs, R. A., Jordan, M. I., Nowlan, S. J., & Hinton, G. E. (1991). Adaptive mixtures of
150 local experts. *Neural Computation*, 3(1), 79–87.
- 151 Liu, Y., Ott, M., Goyal, N., Du, J., Joshi, M., Chen, D., Levy, O., Lewis, M., Zettlemoyer,
152 L., & Stoyanov, V. (2019). RoBERTa: A robustly optimized BERT pretraining approach.
153 *arXiv Preprint arXiv:1907.11692*.
- 154 OpenAI. (2023). GPT-4 technical report. *arXiv Preprint arXiv:2303.08774*.

- 155 Paszke, A., Gross, S., Massa, F., Lerer, A., Bradbury, J., Chanan, G., Killeen, T., Lin,
156 Z., Gimelshein, N., Antiga, L., & others. (2019). PyTorch: An imperative style, high-
157 performance deep learning library. *Advances in Neural Information Processing Systems*,
158 32.
- 159 Pedregosa, F., Varoquaux, G., Gramfort, A., Michel, V., Thirion, B., Grisel, O., Blondel, M.,
160 Prettenhofer, P., Weiss, R., Dubourg, V., & others. (2011). Scikit-learn: Machine learning
161 in python. *Journal of Machine Learning Research*, 12(Oct), 2825–2830.
- 162 Wolf, T., Debut, L., Sanh, V., Chaumond, J., Delangue, C., Moi, A., Cistac, P., Rault, T.,
163 Louf, R., Funtowicz, M., & others. (2019). HuggingFace's transformers: State-of-the-art
164 natural language processing. *arXiv Preprint arXiv:1910.03771*.
- 165 Zadrozny, B., & Elkan, C. (2002). Transforming classifier scores into accurate multiclass
166 probability estimates. *Proceedings of the Eighth ACM SIGKDD International Conference
167 on Knowledge Discovery and Data Mining*, 694–699.
- 168 Zhang, X., Zhao, J., & LeCun, Y. (2015). Character-level convolutional networks for text
169 classification. *Advances in Neural Information Processing Systems*, 28, 649–657.