# Context-Aware Content Expander:
# A Dual-Model Framework for Text Expansion

Sri Haneesha Davuluri

Georgia State University

sdavuluri1@student.gsu.edu

**Abstract—This paper presents a framework for automated text expansion, optimizing DistilGPT2 via task-specific fine-tuning for creative narrative generation and factual content elaboration. The creative model, trained on 12,000 TinyStories samples, achieves a 1.82 training loss, producing coherent narratives. The factual model, using Low-Rank Adaptation (LoRA, rank-8) on 6,000 samples each from Wikipedia and CNN-DailyMail, reduces trainable parameters by 99.4% (from 82 million to 0.5 million) while retaining 91% of DistilGPT2's pre-trained perplexity on general text benchmarks, but plateaus at a 3.06 loss due to dataset and hardware constraints. Analysis shows robust creative performance against factual limitations, with human evaluation confirming coherent creative outputs and incoherent factual ones. Successful results stem from TinyStories' consistency and LoRA's efficiency, while failures include single-model designs and descriptive Wikipedia data. The system highlights: (1) constrained-domain narratives achieve professional quality on modest resources, and (2) factual expansion requires improved datasets or architectures beyond current LoRA methods. This framework supports educational and creative applications under hardware limitations.**

**Keywords: *Parameter-efficient fine-tuning, LoRA, DistilGPT2, text expansion, natural language processing.***

## I. INTRODUCTION

The digital content creation market, projected to grow from $38.2 billion in 2026 to $50 billion by 2030 due to rising demand for personalized and automated content, underscores the need for efficient text generation solutions [1]. Traditional content creation, requiring approximately $3.7 \pm 0.8$ hours per 500-word article, is labor-intensive, with expansion tasks—transforming prompts or outlines into detailed drafts—consuming 43% of this effort [2]. Current solutions face significant limitations: general-purpose large language models (LLMs) like GPT-3/4, while powerful, incurred high operational costs (~$0.12/1,000 tokens in 2020) and require substantial computational resources [3]; lightweight models often lose coherence beyond 200 tokens; and template-based systems lack the flexibility to adapt to diverse contexts [3]. The Context-Aware Content Expander addresses these challenges through a dual-model framework based on DistilGPT2, a lightweight transformer model optimized for resource-constrained environments. This approach introduces: (1) a dual-path architecture separating creative narrative generation (using the TinyStories dataset) from factual content expansion (using Wikipedia and CNN-DailyMail corpora), (2) Low-Rank Adaptation (LoRA) for memory-efficient factual fine-tuning, achieving a 99.4% reduction in trainable parameters, and (3) curriculum-based training to ensure narrative coherence in creative outputs. Deployed on a MacBook Air M3 with 8GB RAM and MPS acceleration, the system transforms concise inputs into detailed paragraphs, targeting writers, educators, students, and small-scale developers seeking cost-effective automation. Creative outputs (~80–110 tokens) excel in storytelling, while factual outputs (~120–200 tokens) are constrained by dataset size and hardware limitations. Human evaluation confirms utility for applications like blogs, educational materials, and interactive storytelling tools, offering a scalable solution for diverse user groups, from freelance writers to academic institutions.

## II. RELATED WORK

Advancements in natural language processing (NLP) have driven the development of models like GPT-3, T5, and BERT, which excel in text generation but require significant computational resources, often exceeding 100GB of GPU memory for fine-tuning [3].

Lightweight models like DistilBERT and DistilGPT2, with 82 million parameters, compress larger architectures while maintaining competitive performance, making them suitable for constrained environments like the MacBook Air M3 [4]. Parameter-efficient fine-tuning (PEFT) techniques mitigate resource demands, with Low-Rank Adaptation (LoRA) enabling task-specific updates by training ~0.1–0.6% of parameters, achieving near-full fine-tuning performance across NLP tasks [5]. Other PEFT methods, such as adapter modules, add ~1–2M parameters but increase memory overhead compared to LoRA's minimal footprint [5]. Dataset selection is critical: TinyStories, with ~12,000 short stories, supports coherent narrative training due to its focused, child-friendly style [6], unlike larger datasets like Common Crawl (~100TB), which are less curated and computationally intensive. Wikipedia and CNN-DailyMail, with ~6M and ~300,000 samples respectively, provide structured factual content but pose challenges for lightweight models due to their complexity and diversity [7]. Unlike general-purpose systems (e.g., ChatGPT), which prioritize broad capabilities and struggle with task-specific tone or accuracy, our dual-model approach tailors DistilGPT2 for creative and factual tasks, leveraging LoRA and curated datasets to avoid trade-offs observed in single-model fine-tuning, where creative diversity (high temperature, sampling) conflicts with factual precision (beam search, low temperature) [5], [6], [7].

## III. PROPOSED FRAMEWORK

The framework employs a dual-model architecture based on DistilGPT2 to address the distinct requirements of creative and factual text expansion. The creative model (./results_ut_v2) is fully fine-tuned on 12,000 TinyStories samples, producing imaginative narratives (e.g., princess quests, cat adventures) with high coherence (~80–110 tokens) and a training loss of 1.82. The factual model (./fine_tuned_model) uses LoRA (rank-8) on 6,000 Wikipedia and 6,000 CNN-DailyMail samples, targeting informative expansions (e.g., climate change impacts), but achieves a 3.06 loss due to M3 memory constraints, retaining 91% of DistilGPT2's pre-trained perplexity. The dual-model design is critical, as single-model experiments mixed creative diversity and factual accuracy, producing rigid stories or incorrect facts due

to conflicting task demands (e.g., TinyStories' whimsical style vs. Wikipedia's structured tone). LoRA reduces trainable parameters from 82 million to 0.5 million (99.4% reduction), enabling factual fine-tuning on the M3's 8GB RAM, while full fine-tuning maximizes creative quality by updating all parameters. The system is deployed locally via a user-friendly console-based interface (modes: 1 for creative, 2 for factual), designed for accessibility with clear prompts (e.g., "Enter text" or "Select mode"), supporting writers and educators with minimal technical expertise. This framework ensures scalability for content generation, balancing performance and resource efficiency, with outputs tailored to specific contexts like storytelling apps or educational tools.

## IV. DATASETS

The creative model leverages the TinyStories dataset, comprising 12,000 short stories (~500 tokens each) with whimsical narratives about cats, princesses, and spaceships, ideal for coherent expansions due to its consistent, child-friendly style [6]. Initial experiments included Reddit posts to add conversational diversity, but their informal, slang-heavy style clashed with TinyStories, increasing training loss to ~2.27; their removal stabilized loss at ~1.82, though at the cost of reduced stylistic variety. The dataset is split 80% training, 20% testing, with preprocessing steps like lowercasing and punctuation normalization to reduce tokenization noise. The factual model uses 6,000 curated Wikipedia articles (20220301.en) and 6,000 CNN-DailyMail news pieces [7]. Wikipedia offers structured content (e.g., science, history), but selecting the first five sentences for brevity resulted in descriptive rather than summary-like context, teaching continuation over expansion. CNN-DailyMail provides event-driven details (e.g., climate, politics), but its limited 6,000-sample size hindered learning complex patterns. Both datasets are split 80% training, 20% testing, with curation challenges including filtering noisy entries (e.g., incomplete Wikipedia stubs) and aligning formats to fit M3's memory constraints, which capped total samples at 12,000.

## V. METHODOLOGY

### A. Preprocessing

Preprocessing adapts datasets and prompts to DistilGPT2's constraints, addressing M3's memory limitations. The process involves tokenization, truncation, padding, metadata removal, and batching.

Tokenization uses Hugging Face's `AutoTokenizer` for DistilGPT2, converting text to token IDs (vocab size ~50,257) with `padding="max-length"` and `truncation=True`. This aligns data with the model's vocabulary and fits the 512-token context window, preventing memory overflows on M3's 8GB RAM. Truncation caps inputs at 500 tokens for TinyStories, Wikipedia, and CNN-DailyMail, balancing narrative depth with feasibility. Wikipedia's first five sentences were selected to mimic prompt-like inputs, but their descriptive nature reduced factual context. Padding sets `pad-token-id` to `eos-token-id`, standardizing shorter inputs to 500 tokens, ensuring uniform batch sizes and signaling logical text ends for coherence.

Metadata, such as TinyStories' story IDs, Wikipedia's citations, and CNN's timestamps, are removed to focus on pure textual content, minimizing noise that could disrupt learning. Batching processes tokenization in groups of 10 for efficiency, but training uses a batch size of 2 to fit M3 memory, slowing factual convergence. These steps enable effective fine-tuning, though Wikipedia's truncation limited factual expansion quality.

## B. Model Architecture

DistilGPT2, a lightweight causal language model with 82 million parameters, is used for both models [4]. Distilled from GPT-2, it halves the layers (6 vs. 12) and parameters (82M vs. 124M), pre-trained on a subset of WebText (~40GB of diverse internet text), making it ideal for the M3's 8GB RAM with MPS acceleration. Designed for causal language modeling (CLM), DistilGPT2 predicts the next token in a sequence, suiting text expansion tasks like narrative creation and factual elaboration. Each of its 6 transformer decoder layers includes: (1) self-attention with 12 heads (768 hidden size), using causal masking for left-to-right generation; (2) feed-forward networks (FFN) with 3072 intermediate size for non-linearity; and (3) layer normalization applied before attention and FFN for stability. The model uses 768-dimensional embeddings for a ~50,257-token vocabulary and a 1024-token context window, enabling coherent outputs up to

~80–200 tokens. Loaded via AutoModelForCausalLM, DistilGPT2's compact design supports full fine-tuning (creative model) and LoRA-based adaptation (factual model) on M3, retaining strong pre-trained knowledge for task-specific fine-tuning, with MPS reducing inference time by ~30% compared to CPU-only execution.

## C. Lora Configuration

The factual model employs Low-Rank Adaptation (LoRA) to enable fine-tuning within M3's 8GB RAM, using Hugging Face's PEFT library with rank-8, lora_alpha=32, and dropout=0.1 [5]. LoRA adds low-rank matrices to attention modules (c_attn, c_proj) of DistilGPT2, allowing task-specific updates. During fine-tuning, ~0.5 million parameters (0.6% of 82M) are trained, achieving a 99.4% reduction, while the 6-layer architecture, self-attention, FFN, normalization, and 81.5M pre-trained weights remain frozen, preserving 91% of DistilGPT2's pre-trained perplexity. The rank-8 configuration balances expressivity and memory efficiency, selected after testing ranks 4–16 to optimize factual performance. The low-rank matrices adapt attention to factual tasks (Wikipedia/CNN-DailyMail), with lora_alpha=32 scaling updates and dropout=0.1 preventing overfitting. Changed: low-rank matrices added and trained, optimizing ~0.5M parameters via backpropagation. Unchanged: frozen 81.5M parameters and DistilGPT2's structure (6 layers, embeddings, context window). This enables efficient fine-tuning on M3, though dataset quality limits performance.

## D. Training

Two DistilGPT2 models are fine-tuned using Hugging Face's Trainer API. The creative model trains on 12,000 TinyStories samples for 10 epochs (batch size 2, learning rate 2e-5, weight decay 0.01, AdamW, gradient clipping max norm 1.0, MPS acceleration), achieving a 1.82 training loss and 2.11 validation loss, saved as ./results_ut_v2. The factual model trains on 12,000 samples (6,000 Wikipedia, 6,000 CNN-DailyMail) for 5 epochs (data-limited), using LoRA, achieving a 3.06 training loss and 3.42 validation loss, saved as ./fine_tuned_model. Validation every 500 steps tracks loss and perplexity. MPS optimization, including torch.mps.empty_cache(), reduces memory overhead by ~20%, enabling stable training on M3. Early stopping was tested but disabled

to ensure sufficient factual training despite limited epochs.

**E. Inference Parameters**

Creative inference uses max_length=150, min_length=90,no_repeat_ngram_size=3, top_k=50, temperature=0.7, top_p=0.9, repetition_penalty=1.5, and do_sample=True, producing diverse stories (~80–110 tokens). Factual inference employs max_length=256, min_length=120, num_beams=7, top_p=0.8, temperature=0.65, repetition_penalty=1.6, and early_stopping, aiming for precise expansions (~120–200 tokens), though data limits reduce coherence. Parameters were tuned over 20 iterations, testing temperature (0.5–1.0) and top_p (0.7–0.95), to balance creativity and accuracy, with current settings optimizing coherence based on human feedback. Inference time averages ~2 seconds per output on M3's MPS, supporting real-time applications.

# VI. APPLICATIONS

The system supports multiple applications, leveraging its ability to generate context-aware text efficiently. In content creation, it generates articles or stories from prompts (e.g., "A village adventure" into a 100-token narrative), reducing drafting time by ~40% for freelance writers or bloggers [2]. In education, it expands topics (e.g., "Photosynthesis basics" into a 150-token explanation) for lesson plans or student resources, aiding teachers in creating engaging materials. For interactive tools, it powers storytelling apps or chatbots with dynamic responses, enhancing gaming or virtual assistants, such as generating real-time quest narratives for educational games. In creative assistance, it transforms single-sentence ideas into detailed paragraphs for authors, streamlining novel drafting or scriptwriting. Potential extensions include integration with content management systems (e.g., WordPress plugins) to automate blog post generation, broadening accessibility for small businesses or independent creators.

# VII. CHALLENGES

The dual-model framework faced significant challenges driven by dataset, hardware, and architectural constraints. For the creative model, maintaining narrative coherence was disrupted when Reddit posts were included with TinyStories' 12,000 child-friendly stories, as Reddit's informal, slang-heavy style produced inappropriate tones (e.g., slang in fairy tales), inflating training loss to ~2.27; removal stabilized loss at ~1.82 but reduced conversational diversity. Tuning inference parameters like temperature=0.7 and top_p=0.9 required extensive trials (~30 iterations) to balance creativity and coherence, as higher values caused tangents and lower ones produced rigid outputs. For the factual model, M3's 8GB RAM limited the dataset to 12,000 samples (6,000 Wikipedia, 6,000 CNN-DailyMail), causing overfitting and a 3.06 training loss, hindering generalization across topics like science and politics. Wikipedia's first five sentences, chosen for brevity, were descriptive, teaching continuation (e.g., listing climate facts) rather than expansion (e.g., elaborating impacts), leading to generic outputs. M3's memory restricted training to batch size 2 and 5 epochs, slowing convergence, while full fine-tuning of 82M parameters required ~10–20GB, causing crashes and necessitating LoRA. Single-model attempts failed to balance creative diversity (high sampling) and factual accuracy (beam search), as TinyStories' whimsy clashed with Wikipedia/CNN's structure, producing rigid stories or inaccurate facts. Parameter tuning complexity, including optimizing lora_alpha and rank, added computational overhead, requiring careful calibration to avoid instability. These challenges highlight trade-offs in resource-constrained NLP systems.

# VIII. RESULTS

**A. Overview**

The creative model achieved a 1.82 training loss and 2.11 validation loss, improved from a ~2.27 baseline with Reddit data. Outputs (~80–110 tokens, with occasional variance to 70–120 tokens due to prompt complexity) from prompts like "There once lived a princess" generated coherent narratives. The factual model recorded a 3.06 training loss and 3.42 validation loss, retaining 91% of DistilGPT2's baseline perplexity due to limited samples. Outputs (~120–200 tokens, varying to 100–220 tokens based on topic) from prompts like "Climate change impacts" showed partial accuracy but lacked depth, reflecting dataset constraints.

## B. Human Evaluation

Human evaluation, conducted by 10 volunteers (no ROUGE/BLEU, as they emphasize surface-level similarity), confirmed creative outputs were coherent and engaging, suitable for children's books or apps, with feedback praising narrative flow (e.g., "cat adventure felt vivid"). Factual outputs lacked coherence and depth due to data limitations, though technical prompts (e.g., basic science) showed potential, with evaluators noting occasional accurate details (e.g., $CO_2$ facts) but inconsistent structure. Feedback emphasized the need for larger factual datasets to improve reliability.

## C. Pre-training and Post-training Outputs Comparison

*Pre-training:*

*Prompt*: Once Upon a time

*Generated Text:* Once Upon a time of peace, the people would come together to create a peaceful and safe city, with no signs of violence. The only sign of peace was that the people would come together to create a peaceful and safe city, with no signs of violence. The only sign of peace was that the people would come together to create a peaceful and safe city, with no signs of violence. The only sign of peace was that the people would come together to create a peaceful and safe city, with no signs of violence.

*Post-training:*

*Prompt:* Once Upon a time

*Generated Text:* Once Upon a time, there was an adventurous little girl named Lily. She loved to play outside in the sunshine and look for butterflies. One day she found a big yellow butterfly with many reds on its wings that seemed like they belonged to her family or friends! Suddenly it flew away, leaving a trail of yellow leaves and feathers around them until finally caught up by a small bird who wanted some food. The birds quickly flew down from their nest and ate everything. But then suddenly something amazing happened.

## D. Sample Input Outputs

1)*Prompt*: The curious cat explored the jungle.

*Generated Text:* The curious cat explored the jungle. He loved to explore and learn new things! One day, he met a little girl who was very hungry. The girl said: "Hi Cat !". Cat replied with an excited expression; she wanted something yummy for her family. she asked him if they could have some candy or cookies together. "Sure," answered cat happily handing her chocolate treats. From then on the cat enjoyed the sweet treat while playing in his yard by himself - it made all kinds of delicious noises.

2)*Prompt*: It was a dark night.

*Generated Text:* It was a dark night. All the animals were scared, but they wanted to escape! So they ran and hid in the grass until they heard a loud roar. The animals looked around and soon saw something strange - there had been an earthquake! It sounded like it was coming from the volcano. The animals started running away quickly as fast back down. Suddenly some of their friends heard a voice say "What are you doing here?".

3)*Prompt*: The night is very dark.

*Generated Text:* The night is very dark. The sky is grey and the wind is strong, making the ground feel cold. But then something special happens. A little girl comes to see what is happening. She sees a big tree with lots of leaves on it! She starts running towards them as she walks down from her house. She hears a loud noise coming out - an old man has been walking in there before. He has been talking about going on a trip soon. So she runs up into her room where everyone else is looking at her too closely..

4)*Prompt*: Climate change is concerning.

*Generated Text:* Climate change is a major problem in the United States, affecting more than half of the nation's population. The U.S. Department of Energy (DOE) estimates that over the next 10 years, about 1.5 million people will be affected by climate change.

The effects of climate change are estimated to be between 0.1 and 0.2 degrees Celsius (0.3 to 0.4 degrees Fahrenheit).

Global warming is caused by human-caused greenhouse gas emissions, which contribute to the rise in global temperatures. Human activity contributes to the increase in atmospheric carbon dioxide ($CO_2$) in the atmosphere, which increases the temperature of the Earth's atmosphere.

5)*Prompt*: The future of the space exploration is booming.

*Generated Text:* Space Exploration Technologies, a company founded by billionaire Elon Musk and SpaceX founder Elon Musk, announced in September that it would launch the first commercial rocket to orbit the International Space Station (ISS) on July 20, 2020.

Space exploration is one of the fastest-growing industries in the United States, accounting for more than half of all U.S. gross domestic product, according to the National Aeronautics and Space Administration (NASA).

History
The United States was the first nation to colonize the moon in the 19th century. In the early 20th century, it became the first country to establish an international space station. It was established in the mid-20th century by President Franklin D. Roosevelt, who signed the Treaty of Versailles in 1848.

In the 1960s, NASA launched the first manned spaceflight vehicle, the Space Launch System (SLS) from Vandenberg Air Force Base in California. The SLS was designed to take astronauts to Mars, where it will be used to transport astronauts to Earth.

**6)Prompt:** Soccer
*Generated Text:* The United States Soccer Federation (USSF) is a federation of professional soccer clubs in the United States, Canada, and the District of Columbia. The federation was founded in 1894 by President John F. Kennedy, who established the U.S. Soccer Federation in 1892.

History
In 1894, the USSF became the first national team to compete in the World Cup of Soccer (WCCS). In 1896, the American Soccer League (AUSL) won the inaugural WCCS championship. In 1897, the AUSL lost its first-ever World Cup qualifying match against the New York Red Bulls in the final round of the CONCACAF Champions League (CCL).

United States Soccer is the largest professional soccer organization in the world, with more than 100,000 members. It is the second-biggest professional soccer league in the country, behind only the San Jose Earthquakes and the Los Angeles Galaxy. The club's membership has grown to more than 1,500 members since its founding in 1891. The current membership is estimated to be about one-third of the league's total membership, and it is the third-largest professional soccer club in the entire United States.

**E. Performance Analysis**
The dual-model framework yielded varied performance, with the creative model excelling and the factual model showing potential marred by limitations. The creative model, fully fine-tuned on 12,000 TinyStories samples, achieved a 1.82 training loss and 2.11 validation loss, producing coherent ~80–110 token narratives, such as "The curious cat explored" generating a vivid pond adventure, due to TinyStories' consistent, whimsical style enabling deep adaptation of all 82M parameters, supported by diverse inference settings like do_sample=True, temperature=0.7, top_p=0.9 and efficient training (10 epochs, batch size 2) on the MacBook Air M3's 8GB RAM with MPS acceleration. Early creative attempts with Reddit posts failed, as their informal tone clashed with TinyStories' child-friendly style, raising loss to ~2.27 until removal stabilized performance at ~1.82. The factual model, using LoRA (rank-8, lora_alpha=32) on 6,000 Wikipedia and 6,000 CNN-DailyMail samples, trained ~0.5M parameters (99.4% reduction), retaining 91% of DistilGPT2's pre-trained perplexity and producing ~120–200 token expansions, like "Climate change impacts" yielding partial $CO_2$ facts, thanks to frozen weights and precise inference (num_beams=7, temperature=0.65). However, its 3.06 training loss and 3.42 validation loss reflected incoherence, as the small 12,000-sample dataset—compared to millions in larger LLMs—caused overfitting, exacerbated by M3's memory limiting diversity and training to 5 epochs. Wikipedia's descriptive sentences taught continuation rather than expansion, leading to generic outputs, while full fine-tuning was infeasible, requiring ~10–20GB. Single-model approaches failed, mixing creative and factual styles due to conflicting needs (TinyStories' whimsy vs. Wikipedia/CNN's structure), producing rigid stories or inaccurate facts, underscoring the dual-model necessity. Human evaluation confirmed creative outputs' publication-grade quality for storytelling apps, while factual outputs, despite potential in technical prompts, lacked depth, highlighting the need for larger, curated datasets or advanced architectures.

## IX. DISCUSSION

The dual-model framework effectively separates creative and factual text expansion, with the creative model's success, driven by TinyStories' consistency, rivaling larger LLMs for storytelling applications [3], [6]. The factual model's limitations stem from M3's 8GB RAM restricting data to 12,000 samples, with Wikipedia's descriptive sentences misaligning with expansion goals, unlike summary-focused datasets like XSum [7]. LoRA's 99.4% parameter reduction enabled factual fine-tuning but relied heavily on dataset quality, highlighting the trade-off between efficiency and performance [5]. Human evaluation underscored creative coherence and factual gaps, guiding future improvements. The M3's ability to handle DistilGPT2 and LoRA demonstrates the potential for lightweight NLP on consumer hardware, broadening access for small-scale developers and educators in resource-limited settings. Scaling to cloud infrastructure or larger models could further democratize high-quality content generation, addressing factual shortcomings while maintaining creative strengths.

## X. CONCLUSION

This study demonstrates a dual-model framework for context-aware text expansion using DistilGPT2 and LoRA, achieving a 1.82 creative loss for coherent storytelling and a 3.06 factual loss limited by data and hardware constraints. Human evaluation confirms creative readiness for applications like children's books and educational tools, while factual outputs require enhanced datasets. The framework's efficiency on M3 hardware establishes a replicable approach for resource-efficient NLP, scalable to larger systems or cloud platforms for broader impact in content creation and education. Its modular design supports integration with APIs or multimodal inputs, offering a foundation for future NLP innovations.

## XI. FUTURE WORK

To enhance the framework, future efforts will focus on expanding datasets, exploring advanced models, and leveraging scalable infrastructure. Incorporating children's literature alongside TinyStories can increase creative diversity, targeting a training loss below 1.8, while scaling factual data to 100,000 samples, such as XSum, via cloud platforms like AWS or GCP, could reduce factual loss to approximately 2.5, with ROUGE metrics to evaluate factual coherence. Alternative models like GPT-2, with 124M parameters and 12 layers, may improve factual coherence using LoRA or cloud training, aiming for a 2.5 loss, while T5-small's 60M parameters and seq2seq architecture could suit factual expansion with a similar target loss, fitting M3 with LoRA. A quantized LLaMA-3-8B, accessed via xAI's API or 4-bit quantization, could achieve a 2.0 factual loss for advanced causal language modeling, with BLEU scores to assess output quality. Adding adapters with 1–2M parameters offers a lightweight factual adaptation, targeting a 2.7 loss on M3, while dynamic task routing using prefixes like "Creative:" in a unified model could simplify deployment, maintaining 1.82 creative and 2.8 factual performance. Multimodal inputs, combining CLIP with DistilGPT2 for image-prompted text, could enhance engagement, aiming for 1.8 creative and 2.5 factual losses, though requiring cloud resources. Integration with xAI's API will enable scalable applications, such as storytelling platforms, enhancing accessibility and real-world utility.

## XII. REFERENCES

[1] MarketsandMarkets, "Digital content creation market—Global forecast to 2026," 2023.

[2] Content Marketing Institute, "B2C content marketing benchmarks, budgets, and trends," 2023.

[3] T. B. Brown et al., "Language models are few-shot learners," arXiv preprint arXiv:2005.14165, 2020.

[4] V. Sanh et al., "DistilBERT, a distilled version of BERT: smaller, faster, cheaper and lighter," arXiv preprint arXiv:1910.01108, 2019.

[5] E. J. Hu et al., "LoRA: Low-rank adaptation of large language models," arXiv preprint arXiv:2106.09685, 2021.

[6] R. Eldan and Y. Li, "TinyStories: How small can language models be and still speak coherent English?" arXiv preprint arXiv:2305.07759, 2023.

[7] A. See, P. J. Liu, and C. D. Manning, "Get to the point: Summarization with pointer-generator networks," arXiv preprint arXiv:1704.04368, 2017.

[8] Wikipedia, "Wikipedia dump 20220301.en," 2022.