

OpenMic: A Multi-Agent-Based Stand-Up Comedy Generation System

Yuyang Wu¹ Hanzhong Cao¹ Jianhao Chen¹ Yufei Li¹

Abstract

Chinese stand-up comedy generation goes beyond plain text generation, requiring culturally grounded humor, precise timing, stage-performance cues, and implicit multi-step reasoning. Moreover, commonly used Chinese humor datasets are often better suited for humor understanding and evaluation than for long-form stand-up generation, making direct supervision misaligned with the target task. To address these challenges, we present **OpenMic**, an end-to-end multi-agent system built on AutoGen that transforms a user-provided life topic into a 3–5 minute Chinese stand-up performance and further produces a narrated comedy video. OpenMic orchestrates multiple specialized agents in a multi-round iterative loop—planning to jointly optimize humor, timing, and performability. To mitigate the dataset–task mismatch, we augment generation with retrieval-augmented generation (RAG) for material grounding and idea expansion, and we fine-tune a dedicated JokeWriter to better internalize stand-up-specific setup–punchline structures and long-range callbacks.

1. Introduction

Artificial intelligence has made rapid progress in creative content generation, spanning text, music, and visual art. Yet performative creativity remains notably harder: stand-up comedy is not just “good writing,” but a tightly choreographed sequence of linguistic craft, temporal control, and social-context awareness. This gap is reflected even in industry practice—despite the scale and capability of modern foundation models, major labs rarely report standardized “humor ability,” largely because humor evaluation itself is intrinsically difficult: what counts as funny is subjective, culturally grounded, context-dependent, and highly sensitive to delivery and timing.

¹School of Electronics Engineering and Computer Science, Peking University. Correspondence to: Yuyang Wu <wuyuyang@stu.pku.edu.cn>.

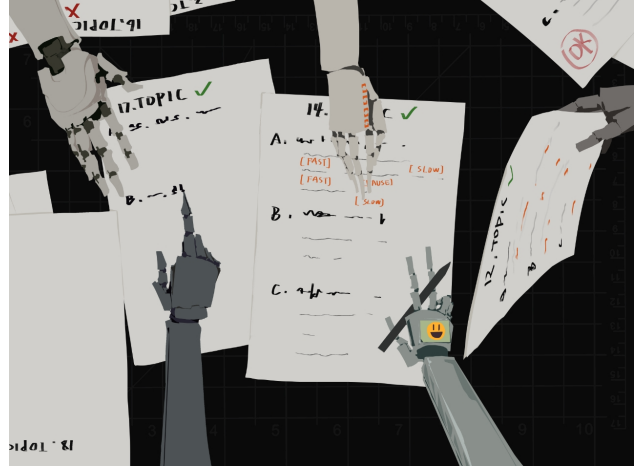


Figure 1. Multi-agent collaborative pipeline for Chinese stand-up comedy generation. Specialized agents iteratively decompose ideation, retrieval, joke writing.

Recent research therefore tends to emphasize humor understanding and evaluation rather than full humor generation, because the former is easier to define and benchmark. In Chinese, this tilt is particularly visible: datasets such as CFunSet (Yu et al., 2025) provide rich resources for analyzing and probing humor, with tasks including (i) humor cause analysis, (ii) crosstalk “straight-man” response, (iii) humor binary classification, (iv) keyword-based joke generation, (v) topic-conditioned joke generation, and (vi) joke continuation. While these tasks are valuable for modeling comedic signals and building evaluators, they do not directly match the target of long-form Chinese stand-up: a 3–5 minute performance requires coherent comedic arcs, delayed punchlines, callbacks, and stage-ready phrasing—properties that are under-specified by short-form supervision and are hard to learn from understanding-centric labels alone.

Meanwhile, generation remains difficult even with strong general-purpose models. As our preliminary comparison in Fig. 7 and 8 suggests, a strong general model (e.g., GPT-5.2) can drift into didactic or “preachy” narration when asked to produce stand-up, while another strong Chinese model (e.g., DeepSeek) may produce jokes that are sparse and uneven in quality. These failures are not simply stylistic; they reveal missing control over (1) comedic structure (setup–punchline delay, misdirection, callback), (2) timing (pauses, emphasis,

rhythm), and (3) performability (spoken language and stage cues). In other words, humor is not equivalent to fluent text, and Chinese stand-up further amplifies the challenge through heavier reliance on shared social context, colloquial delivery, and timing-sensitive audience expectation management.

To address these issues, we build on the intuition that stand-up generation is closer to a production pipeline than a single-shot completion: it requires planning, audience adaptation, writing, coaching, and critique—each with different objectives and failure modes. We therefore propose OpenMic, a multi-agent system implemented with AutoGen (Wu et al., 2024), where specialized agents collaborate in a multi-round iterative loop to refine content toward both comedic quality and stage readiness. To bridge the dataset–task mismatch, we incorporate retrieval-augmented generation (RAG) to ground writing in diverse comedic materials and to expand topic-specific angles, and we fine-tune a dedicated JokeWriter to better internalize stand-up-oriented structures beyond what understanding-focused datasets naturally provide. Finally, OpenMic outputs not only a script but a structured performance representation (e.g., pauses, applause beats, emphasis) that can be rendered into an end-to-end video.

- We implement an end-to-end multi-agent Chinese stand-up comedy generation system based on AutoGen, from user topic input to a stage-ready performance.
- We introduce RAG-based material retrieval to enrich content grounding and alleviate sparsity/mismatch of stand-up supervision.
- We design a multi-round iterative self-improve workflow to improve comedic structure, timing, and performability.
- We fine-tune a dedicated JokeWriter to better capture setup–punchline delay, callback patterns, and spoken-stage style.
- We propose a structured performance script interface (pauses, applause, emphasis, etc.) and a pipeline that converts it into narrated comedy video output.

2. Related Works

The field of computational humor has long been considered an "AI-complete" problem because it requires a deep understanding of semantics, pragmatics, and social context. (Kim & Chilton, 2025a) Traditional research focused on Incongruity Theory, which posits that humor arises from the sudden resolution of a mismatch between expectations and reality. (Chen et al., 2024; Loakman et al., 2025) In the context of crosstalk and talkshows, this is manifested

as the "set-up and punchline" logic (or Baofu in Chinese crosstalk). Early attempts at humor generation were often template-based and lacked the creative "logic jump" required for effective comedy. (Kim & Chilton, 2025b) Recent work such as "Humor Mechanics: Advancing Humor Generation with Multistep Reasoning" has shifted the focus toward reconstructing these mechanics through data-driven policies. (Tikhonov & Shtykovskiy, 2024) They demonstrate that humor is not merely a linguistic byproduct but a result of multistep reasoning where the model must distill humor principles—such as wordplay and unexpected twists—from existing datasets to generate novel content rather than just acting as a "stochastic parrot."

Our technical framework draws from three rapidly evolving areas of NLP. First, while traditional Retrieval-Augmented Generation (RAG) was primarily used for fact-checking, it has recently been adapted for creative tasks to inject cultural "memes" and specific comedic styles. (Sanmartin, 2024) Current trends favor Hybrid Adaptation (similar to Retrieval-Augmented Fine-Tuning or RAFT), which balances the static domain expertise of the model with dynamic, external context. (Balaguer et al., 2024) Second, the development of Parameter-Efficient Fine-Tuning (PEFT) has moved from LoRA to QLoRA, allowing for the specialization of large language models (LLMs) on high-quality comedic scripts without the prohibitive cost of full retraining. (Dettmers et al., 2023) Finally, our architecture utilizes a Multi-Agent System (MAS) to mimic human collaborative creativity. We build upon the foundation of works like HoLLMwood which assigns LLMs to specialized roles such as "Writer," "Editor," and "Actor" to improve narrative coherence. (Chen et al., 2024) By following the modular design principles outlined in recent MAS surveys, we create a specialized pipeline where different agents handle distinct stages of the crosstalk generation process.

The current landscape of humor generation is increasingly focusing on multi-dimensional evaluation and cultural specificity. For instance, this paper (Sakabe et al., 2025) reveals that while modern LLMs can match low-to-mid tier human performance in improvisational Japanese comedy, they often prioritize "Novelty" over "Empathy," leading to a divergence in what machines and humans perceive as funny. Similarly, Guo et al. (2023) highlighted the gap in LLM performance for Chinese crosstalk, where models struggle with the rhythmic cadence and the specific structural requirements of the medium. (Wang et al., 2022) Our work contributes to this evolving field by combining a novel multi-LLM agent system with a RAG-based context injector and an agent-specific fine-tuning strategy. By specifically training agents to play individual roles (e.g., the "JokeWriter" agent), we explore whether the synergy of specialized roles and retrieved comedic materials can overcome the empathy-novelty gap identified in recent benchmarks.

Q: Why was the cookie sad?
A: Because his mom was a wafer long!

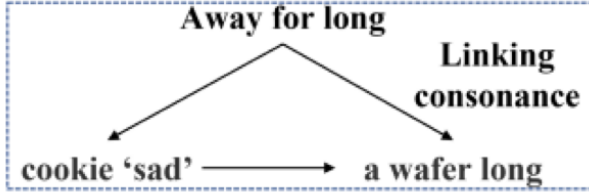


Figure 2. Backdoor Criterion

Q: What does a clock do when it's hungry?
A: It goes back for seconds.

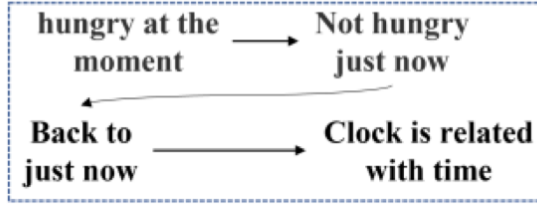


Figure 3. Frontdoor Criterion

3. Humor and Cognitive Reasoning

Humor is often treated as a stylistic property of language, yet many jokes are better understood as *reasoning processes* that manipulate an audience’s expectations over time. A stand-up punchline rarely succeeds by lexical novelty alone; rather, it relies on a latent chain of inferences that (i) builds a plausible interpretation, (ii) introduces a hidden connection, and (iii) triggers a rapid “re-interpretation” that resolves the incongruity. This view is especially important for Chinese stand-up comedy, where effective jokes frequently depend on shared cultural context, implicit premises, and tightly controlled information release. In this section, we frame humor as a form of structured cognitive reasoning and use two illustrative logic diagrams (Fig. 2 and Fig. 3) to highlight distinct inference patterns that commonly appear in jokes.

3.1. Humor as a Form of Cognitive Incongruity

Classic humor theories emphasize *incongruity*: a joke sets up an expectation and then violates it in a way that still permits a coherent resolution. Under an incongruity–resolution perspective, “being funny” is not merely about generating surprising words, but about creating a *controlled mismatch* between the audience’s predicted continuation and the eventual reinterpretation that makes the punchline sensible. Concretely, the setup implicitly constructs a mental model of the situation; the punchline either flips a key assumption or

reveals a hidden linkage that forces the audience to update that model. The comedic effect arises from the *contrast* between the initial expectation and the revised interpretation, as well as the speed and clarity with which the resolution becomes apparent.

3.2. Humor Requires Multi-Step Reasoning

Many jokes embed a multi-step inference chain rather than a single-step association. Fig. 2 provides an example we refer to as a *backdoor-style* structure: the question entity E_Q and answer entity E_A are not directly connected by surface meaning, but are linked through an intermediate bridging entity E_Z (often a homophone, pun, or shared attribute). In the illustrated joke (“Why was the cookie sad?”), the surface reading suggests an emotional explanation; the resolution depends on mapping to the phonetic/lexical bridge (e.g., “away for long” \leftrightarrow “a wafer long”), which then retroactively makes the punchline interpretable. Here, E_Z acts as a hidden connector that is easy to miss unless one actively searches for alternative interpretations.

In contrast, Fig. 3 illustrates a *frontdoor-style* multi-hop reasoning pattern: the setup encourages the audience to traverse intermediate thoughts explicitly before arriving at the punchline. In the example (“What does a clock do when it’s hungry?” \rightarrow “It goes back for seconds.”), the humor hinges on composing several simple steps: “hungry” evokes a desire for food; “seconds” can mean a second helping; and “clock” relates to time, enabling the wordplay “goes back for seconds.” Compared with the backdoor-style pun, the intermediate entity E_Z in this case is not merely a hidden phonetic bridge but a *conceptual stepping stone* that the listener can traverse through associative and compositional reasoning.

These two patterns are common in stand-up: (i) **delayed punchlines** resemble multi-step inference with intentionally withheld bridges, and (ii) **callbacks** resemble long-range reasoning where an earlier premise is reactivated under a new interpretation. As a result, humor quality depends not only on what is said, but on *when* the crucial bridge is revealed and how reliably the audience can reconstruct the implicit reasoning path.

3.3. Why LLMs Struggle with Humor

Despite strong general language ability, LLMs frequently underperform on humor because fluent continuation does not guarantee the *cognitive surprise* required for a joke. First, models tend to collapse the inference process: they may reveal the bridge too early, explain the joke explicitly, or smooth over ambiguity—all of which reduce comedic tension. Second, many jokes require maintaining two competing interpretations until the punchline; this demands deliberate control of uncertainty and information release, whereas

next-token prediction often favors a single dominant continuation. Third, humor is highly sensitive to pragmatic constraints (social norms, persona, cultural presuppositions), so even logically consistent outputs may fail to land as funny if the implied premises are unnatural for the target audience.

The reasoning structures in Fig. 2–3 also expose a practical issue: the model must *search* over potential bridges E_Z (phonetic, semantic, or contextual) and then *stage* the reveal at the right moment. Without explicit mechanisms for planning, critique, and timing control, single-shot generation often produces either (i) coherent but unfunny narration, or (ii) isolated one-liners that lack buildup, callbacks, and performance rhythm.

3.4. Implication: Humor as Structured Reasoning

Viewing humor as structured reasoning suggests that effective stand-up generation is closer to a pipeline of **planning**, **verification**, and **execution** than to free-form text generation. Planning selects comedic angles and determines which bridge E_Z to hide or surface and when; verification checks whether the reasoning path is reconstructible and whether the punchline resolves the incongruity; execution adds stage-performance cues (pauses, emphasis, and rhythm) that regulate information release. This framing motivates our system design: instead of asking a single model to simultaneously invent content, manage audience adaptation, and control delivery, we decompose the process into specialized roles and iterative refinement so that the final performance preserves both the reasoning structure and the timing-sensitive comedic payoff.

4. Methodology

4.1. Multi-Agent System Design

Why Multi-Agent for Stand-Up Comedy? As discussed in Sec. 3.1, humor is tightly coupled with *structured reasoning* and *timing-aware information release*. In practice, Chinese stand-up generation involves several competing objectives that are hard to satisfy in a single pass: (i) **content planning** (selecting angles, building setups, placing callbacks), (ii) **audience adaptation** (persona, taboo avoidance, cultural priors), (iii) **performability** (spoken-style wording, rhythm, pauses, emphasis), and (iv) **quality assurance** (coherence, novelty, safety, “laugh potential”). A single-agent LLM prompt tends to entangle these goals and often collapses the latent reasoning structure (e.g., explaining the joke too early) or drifts into generic narration.

We therefore adopt a multi-agent design that decomposes the pipeline into specialized roles with explicit responsibilities. This separation yields two practical advantages: (1) **controllability**—each agent optimizes a well-defined sub-objective with dedicated constraints; and (2) **iterative**

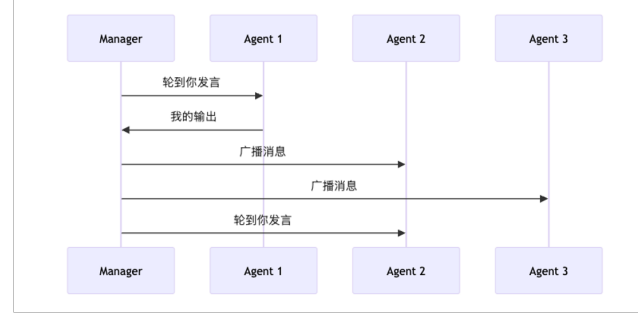


Figure 4. Mechanism of AutoGen

refinement—failures can be localized and corrected (e.g., rewrite a weak punchline without redoing audience profiling), which is essential for timing-sensitive comedic arcs.

4.2. AutoGen as the Orchestration Backbone

We implement OpenMic on top of AutoGen’s group-chat abstraction. Conceptually, AutoGen coordinates a set of conversable agents via a manager that (i) decides whose turn it is to speak, (ii) collects the agent’s output, and (iii) broadcasts relevant messages to other agents for subsequent turns (Fig. 4). This “turn-taking + broadcast” mechanism is a natural fit for creative collaboration: agents can operate asynchronously in intent (each has its own rubric), yet remain synchronized through shared context.

In our implementation, we use a **GroupChatManager** to enforce an ordered protocol and to prevent uncontrolled multi-agent chatter. Each agent is instantiated as a **ConversableAgent** with a role-specific system prompt, input/output schema, and constraints. The manager schedules agents following our workflow (Sec. 4.6) and handles termination conditions (either `PASS` from the quality controller or a maximum iteration budget).

4.3. Blackboard-Centric Coordination

Beyond message passing, OpenMic employs a **blackboard** to maintain structured shared state (Fig. 5). The blackboard stores intermediate artifacts that must persist across turns and iterations, including:

- **Audience profile:** persona, preferences, taboo list, acceptable language register;
- **Topic expansion:** subtopics, personal anecdotes angles, candidate premises;
- **Draft script:** current version of the stand-up text with section boundaries;
- **Performance markup:** a structured DSL with pauses,

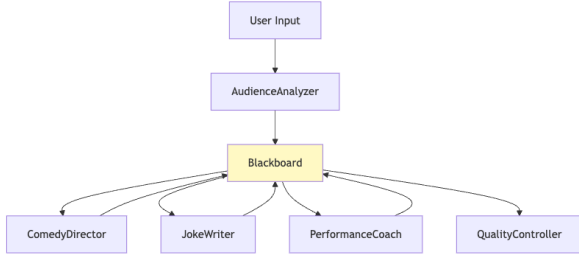


Figure 5. Multi-agent Structure.

emphasis, applause beats;

- **Critique & action items:** concrete revision instructions and failure reasons.

This design prevents critical information from being lost in long conversational context and makes the iteration loop more deterministic: each agent reads from and writes to designated blackboard fields, rather than relying solely on implicit conversational memory.

4.4. Agent Roles and Interfaces

OpenMic consists of five core agents (Fig. 5), each with a narrowly-scoped responsibility and a typed output that is written to the blackboard:

AudienceAnalyzer (audience modeling). Given the user topic and optional style constraints, the AudienceAnalyzer produces an audience-facing **persona card** and a **taboo/avoid list** (e.g., sensitive references, overly offensive wording) to ensure cultural and situational appropriateness.

ComedyDirector (high-level planning). The ComedyDirector decomposes the topic into a set of **subtopics** and a **comedic structure plan** (e.g., opening hook, 2–3 bits, a callback, closing tag). The output is an outline with explicit comedic intent: where tension is built, where bridges are revealed, and where callbacks should land.

JokeWriter (script drafting). Conditioned on the outline and audience profile, the JokeWriter produces a **complete draft script**. We instruct it to maintain spoken Chinese, enforce setup–punchline delays, and preserve long-range dependencies (callbacks) rather than generating disconnected one-liners.

PerformanceCoach (delivery & markup). The PerformanceCoach transforms the draft into a **performance-ready script** by adding a structured DSL annotation, including pauses (e.g., pre-punchline pause), emphasis, pace changes, filler words, and optional applause/laughter cues. This

bridges text generation with downstream audio/video rendering.

QualityController (evaluation & gating). The QualityController acts as a critic and gatekeeper. It evaluates coherence, comedic payoff, timing realism, and audience fit, then outputs either **PASS** or **REVISION** with actionable edits. This turns subjective humor quality into an operational criterion for iteration.

4.5. Hierarchical Multi-Agent RAG with Information Isolation

Our RAG framework is designed to bridge the gap between simple semantic retrieval and genuine creative transformation. Rather than relying on a traditional single-step retrieval process, we implemented a **triadic inner-conversation architecture** that utilizes one retrieval engine alongside two specialized LLM agents. This system is governed by a custom protocol to ensure that the massive volume of data required for candidate selection does not clutter the primary workflow’s context window.

Dataset Composition and Post-Processing To ensure stylistic consistency, our retrieval corpus combines two primary sources. The first is a collection of short-form setups and punchlines sourced from the CFUN repository (Yu et al., 2025). The second is a **Crosstalk-to-Talkshow Pipeline** where we took traditional crosstalk scripts and pushed them through an LLM-driven refinement stage. During this process, we performed anonymization by removing specific performer names and executed a stylistic conversion. This turned dialogue-heavy routines into narrative-driven talkshow observations, moving away from the classic “teasing and reacting” dynamic to a more modern first-person perspective.

The Triadic Inner-Conversation Workflow Standard semantic matching often prioritizes factual similarity over comedic value. To fix this, we formalize the RAG process as a sequence of three specialized operations: retrieval, scoring, and refinement.

Let q represent the user topic query and \mathcal{D} the integrated comedic corpus. We define $E(\cdot)$ as the embedding function that maps text to a high-dimensional vector space. The process is defined as follows:

1. **Semantic Retrieval:** The RAG Retriever identifies a set of raw candidates \mathcal{C} by calculating the cosine similarity between the query and document embeddings:

$$\mathcal{C} = \text{top-}k_1 \{d \in \mathcal{D} \mid \text{sim}(E(q), E(d))\}$$

$$\text{where } \text{sim}(\mathbf{u}, \mathbf{v}) = \frac{\mathbf{u} \cdot \mathbf{v}}{\|\mathbf{u}\| \|\mathbf{v}\|}.$$

2. **LLM Candidate Scoring:** The **LLM Candidate Scorer**

(Agent 1) acts as a non-linear semantic filter. It evaluates the comedic potential P of each candidate $c \in \mathcal{C}$ based on latent features like incongruity and relevance, selecting a subset \mathcal{S} of high-potency jokes:

$$\mathcal{S} = \{c \in \mathcal{C} \mid f_{\text{scorer}}(c, q) > \tau\}$$

where f_{scorer} is the agent’s internal evaluation function and τ is the quality threshold for the top- k_2 selection.

3. **LLM Punchline Refinement:** Finally, the **LLM Punchline Selector** (Agent 2) performs the creative transformation \mathcal{T} . Instead of passing the full text, it distills the selected jokes into a set of writing materials \mathcal{M} :

$$\mathcal{M} = \bigcup_{s \in \mathcal{S}} \mathcal{T}_{\text{selector}}(s)$$

This ensures the JokeWriter receives a distilled set of high-potency building blocks \mathcal{M} rather than a wall of raw, unorganized text, significantly reducing context noise while maximizing creative signal.

The ”Secret Blackboard” and Context Management A key technical feature of our architecture is the **Secret Blackboard**. During the inner-conversation between the RAG engine and the retrieval agents, thousands of tokens of raw material are processed simultaneously. Storing this in the main global blackboard would quickly exceed the context limits of agents further down the line, such as the Performance Coach. To solve this, the Secret Blackboard acts as a private memory buffer. It only releases the final refined punchlines to the JokeWriter, effectively hiding the noisy retrieval process from the rest of the chain and maintaining a high signal-to-noise ratio across the entire system.

4.6. Multi-Round Refinement

If the QualityController returns `REVISION`, the system re-enters the loop by routing feedback back to the JokeWriter (and optionally the PerformanceCoach) until either `PASS` is obtained or a maximum number of rounds is reached. This multi-round loop is crucial for stand-up: jokes often fail due to localized issues (weak punchline, premature explanation, missing callback trigger, unnatural pause placement) that are best fixed through targeted rewrites rather than regenerating everything from scratch.

Dual-Dimension Quality Assessment Unlike monolithic QA systems, our QualityController performs **dual-dimension evaluation** $Q_r = (Q_r^R, Q_r^W)$ to separately assess retrieval quality and writer quality.

RAG Dimension (Q_r^R): Evaluates retrieved joke material quality—humor potential, topic relevance, and diversity. When $Q_r^R = 0$, the QA outputs:

- \mathbf{k}^* : refined keywords
- \mathcal{E}_r : joke IDs to exclude in next retrieval
- f_r^R : specific feedback per joke

Writer Dimension (Q_r^W): Evaluates script organization via three checks:

$$Q_r^W = \mathbb{K}[\text{struct}] \wedge \mathbb{K}[\text{safe}] \wedge \mathbb{K}[\text{length}] \quad (1)$$

Failed checks trigger *rewrite directives* \mathbf{d}_r (e.g., “callback missing for setup in line 3”).

Targeted Refinement Routing The dual evaluation enables surgical fixes:

Case 1: RAG fails, Writer succeeds ($Q_r^R = 0, Q_r^W = 1$):

$$\text{Re-retrieve: } \mathcal{D}_{r+1} \leftarrow \text{RAG}(\mathbf{k}^*, \mathcal{E}_r) \quad (2)$$

Case 2: Writer fails, RAG succeeds ($Q_r^R = 1, Q_r^W = 0$):

$$\text{Rewrite: } s_{r+1} \leftarrow \text{Writer}(\mathcal{D}_r, \mathbf{d}_r) \quad (3)$$

Case 3: Both fail ($Q_r^R = 0, Q_r^W = 0$):

$$\mathcal{D}_{r+1}, s_{r+1} \leftarrow \text{RAG}(\mathbf{k}^*, \mathcal{E}_r) + \text{Writer}(\cdot, \mathbf{d}_r) \quad (4)$$

Termination occurs when $Q_r^R \wedge Q_r^W = 1$ or $r \geq R_{\max}$.

Context-Aware Memory Each writer receives structured feedback:

$$\mathcal{C}_r = \{s_{r-1}, \mathcal{D}_{r-1}, \mathbf{d}_{r-1}, Q_{r-1}^R, \mathcal{P}_{r-1}\} \quad (5)$$

where \mathcal{P}_{r-1} are preserved joke IDs. This prevents agents from “forgetting” prior decisions across rounds.

Empirically, Case 3 occurs in $\sim 30\%$ of round-1 attempts but drops to $< 5\%$ by round 3, indicating rapid convergence.

4.7. Domain-Specific Adaptation via QLoRA

To bridge the gap between general-language capabilities and specialized comedic timing, we employ **Quantized Low-Rank Adaptation (QLoRA)**. This approach allows for the fine-tuning of large-scale models by injecting trainable low-rank matrices into the frozen, 4-bit quantized base model. For a weight matrix $W_0 \in \mathbb{R}^{d \times k}$, the forward pass is modified as:

$$h = W_0 x + \Delta W x = W_0 x + B A x$$

where $B \in \mathbb{R}^{d \times r}$ and $A \in \mathbb{R}^{r \times k}$ are the low-rank adapters with rank $r \ll \min(d, k)$. We specifically target all **linear projections** within the transformer blocks to maximize

Table 1. Evaluation Results across Different Temperature Settings

Configuration	Persona	Reactivity	Humor	Narrative	Coherence
JW+Tem0.1	82.5	88.0	95.5	93.0	97.5
JW+Tem0.3	88.5	15.0	96.0	93.0	97.5
JW+Tem0.5	82.5	15.0	96.5	93.0	98.5
JW+Tem0.7	92.5	45.0	96.0	98.5	97.0
JW+Tem0.9	85.0	25.0	92.0	94.0	96.0

the model’s stylistic plasticity. Furthermore, to ensure the model focuses exclusively on comedic delivery, we utilize a **completion-only loss strategy**, calculating gradients only on the generated punchlines rather than the instruction prompts.

5. Experiments

5.1. LLM-as-a-Judge Evaluation Framework

We implemented a rigorous “LLM-as-a-Judge” mechanism to quantify the quality of the generated talk show scripts. Unlike standard NLP metrics (such as BLEU or ROUGE), which often fail to capture the semantic nuance and comedic timing of creative writing, we utilized a senior executive producer persona—powered by the Grok-4-1-fast-reasoning model—to conduct a multi-dimensional scoring analysis. The evaluation is governed by a Pydantic-enforced schema, ensuring that every assessment is structured across five critical dimensions:

- **Persona Fidelity (30%):** The distinctiveness and consistency of the characters’ voices.
- **Humor Mechanics (25%):** The density and structural quality of setup-punchline sequences.
- **Interactive Reactivity (20%):** The degree of “improvisational” riffing and response to previous turns.
- **Contextual Coherence (15%):** The logical consistency and effective use of callbacks.
- **Narrative Arc (10%):** The rhythmic flow from introduction to climax and resolution.

The final weighted score S_{total} is calculated as:

$$S_{total} = 0.30P + 0.25H + 0.20R + 0.15C + 0.10N$$

where P, H, R, C, N represent the scores for Persona, Humor, Reactivity, Coherence, and Narrative respectively.

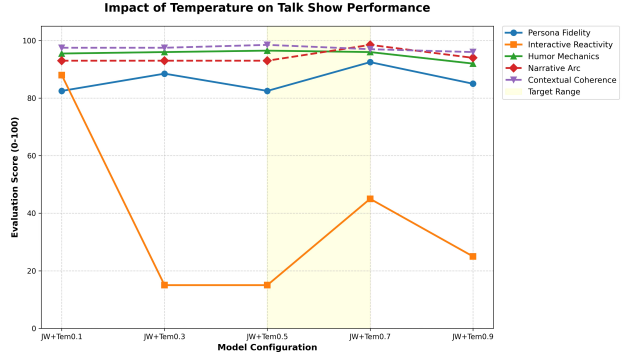


Figure 6. visualization of generation scores under different temperature paramters

5.2. Influence of Temperature on Generative Performance

While initial metrics exhibited a high degree of variance, they provided critical insights into the relationship between sampling temperature and the efficacy of our RAG-enhanced retrieval. We observed that the system’s performance is not a linear function of stochasticity, but rather a delicate balance between **contextual grounding** and creative divergence.

Our analysis indicates that **low temperature settings** (e.g., $T = 0.1$) combined with the **large joke corpus retrieved via RAG** yield the most superior results. As shown in 10, which compares two specific generative examples, lower temperatures allow the model to maintain a high “focus” on the specific comedic building blocks provided by the RAG inner-conversation. In this regime, the **Contextual Coherence** (97.5) and **Interactive Reactivity** (88.0) are maximized, as the model accurately maps the retrieved punchlines onto the target persona without drifting into irrelevant hallucinations.

Conversely, at **higher temperatures** (e.g., $T \geq 0.7$), the model begins to lose its grip on the retrieved context. While this occasionally results in a spike in **Narrative Arc** (98.5 at $T = 0.7$) as the model explores more varied sentence structures, it frequently compromises **Reactivity**. Qualitative review suggests that at high temperatures, the JokeWriter often ignores the specific “setup” provided by the RAG selector in favor of generic, less interesting tropes.

Ultimately, we conclude that the optimal configuration for generative crosstalk lies in **minimizing entropy to maximize retrieval signal**. By utilizing a low temperature, we ensure that the fine-tuned model acts as a precision instrument that “assembles” the retrieved comedic materials into a cohesive script, rather than attempting to hallucinate humor without sufficient grounding. This reinforces the value of our RAG-centric approach: the “creativity” is supplied by the diversity of the corpus, while the “logic” is preserved by the

constrained sampling.

5.3. Finetuning implementation Details and Hyperparameters

Our fine-tuning experiments were conducted on a single GPU using the `trl` and `peft` libraries. The training corpus consists of the LLM-processed Talkshow dataset, formatted using the Qwen-2.5 chat template.

Fine-tuning Configuration: We utilized the `LoraConfig` to target a comprehensive set of modules, including `q_proj`, `k_proj`, `v_proj`, `o_proj`, and the MLP layers (`gate`, `up`, `down_proj`). The detailed hyperparameter settings are summarized in Table 2.

Table 2. Hyperparameters for Comedy-Specialized QLoRA Fine-tuning

Hyperparameter	Value
Base Model	Qwen-2.5-3B-Instruct
Quantization	4-bit NF4 (NormalFloat)
LoRA Rank (r)	16
LoRA Alpha (α)	32
LoRA Dropout	0.05
Learning Rate	2×10^{-4}
Optimizer	Paged AdamW 32-bit
Batch Size (Per Device)	2
Gradient Accumulation	4
Training Epochs	1
Compute Precision	BF16 (or FP16)

Completion-Only Training: To prevent the model from overfitting on the instruction syntax, we implemented a **ManualCompletionCollator**. By defining the response template as "`<|im_start|>assistant\n`", the trainer effectively masks the prompt tokens during loss calculation. This ensures the negative log-likelihood loss \mathcal{L} is computed only on the tokens y_i belonging to the assistant’s response:

$$\mathcal{L} = - \sum_{i \in \text{Response}} \log P(y_i \mid y_{<i}, \text{Prompt})$$

The final model was deployed via a vLLM endpoint with LoRA support enabled, allowing for high-throughput inference during the multi-agent execution phase.

5.4. Demo: End-to-End Stand-Up Generation

Goal. We present an end-to-end demo to illustrate how our multi-agent pipeline generates a Chinese stand-up monologue from a high-level prompt, highlighting (i) structured setup–punchline planning, (ii) iterative critique and rewriting, and (iii) callback triggering across the script.

Setup. We run the system for 3 iterations producing intermediate artifacts including *Quality Evaluations*.

5.5. Downstream Application: End-to-End Video Synthesis

To further demonstrate the practical utility of OPENMIC, we extend the pipeline to a multi-modal application stage. This stage verifies that the structured performance scripts, enriched with behavioral cues, can be seamlessly executed by external rendering engines to produce broadcast-ready content.

Implementation Workflow: The video synthesis process acts as a specialized consumer of the *PerformanceCoach*’s output. We implement a middleware that parses the embedded DSL markers—such as `[pause]`, `[emphasis]`, and `[applause]`—to construct a synchronized temporal timeline. By invoking RESTful APIs from high-fidelity digital human platforms (e.g., Kling AI), the system maps the synthesized audio onto a 3D-animated avatar. The synchronization logic ensures that the avatar’s micro-expressions, such as eyebrow movements during a setup and a smirk during a punchline, are aligned with the comedic rhythm defined in the script.

Key Technical Challenges and Observations:

- **Temporal Consistency:** The use of structured markers prevents the common “robotic delivery” seen in standard Text-to-Speech (T2S) systems. By explicitly injecting silence durations and speech rate variations based on the DSL, we preserve the timing-sensitive nature of Chinese stand-up comedy.
- **Cross-Modal Stylistic Alignment:** The visual persona, including stage background illumination and character attire, is dynamically selected to match the *AudienceAnalyzer*’s persona card. This ensures a coherent comedic atmosphere where the visual environment reinforces the linguistic tone.
- **Performance Fidelity:** Our pipeline automates the generation of a 3–5 minute narrated video from a single topic prompt. This end-to-end capability demonstrates the robustness of OPENMIC not only as a writing assistant but as a comprehensive production tool for digital entertainment.

The integration of video synthesis completes the generative loop, providing a tangible interface for evaluating the performability of the generated humor in a real-world setting.

Acknowledgments

This project is a course final assignment for the CoRE course. It was developed by the group “King of Comedy”.

References

- Balaguer, A., Benara, V., de Freitas Cunha, R. L., de M. Estevão Filho, R., Hendry, T., Holstein, D., Marsman, J., Mecklenburg, N., Malvar, S., Nunes, L. O., Padilha, R., Sharp, M., Silva, B., Sharma, S., Aski, V., and Chandra, R. Rag vs fine-tuning: Pipelines, tradeoffs, and a case study on agriculture, 2024. URL <https://arxiv.org/abs/2401.08406>.
- Chen, J., Zhu, X., Yang, C., Shi, C., Xi, Y., Zhang, Y., Wang, J., Pu, J., Zhang, R., Yang, Y., and Feng, T. Hollmwood: Unleashing the creativity of large language models in screenwriting via role playing, 2024. URL <https://arxiv.org/abs/2406.11683>.
- Dettmers, T., Pagnoni, A., Holtzman, A., and Zettlemoyer, L. Qlora: Efficient finetuning of quantized llms, 2023. URL <https://arxiv.org/abs/2305.14314>.
- Kim, S. and Chilton, L. *AI Humor Generation: Cognitive, Social and Creative Skills for Effective Humor*. 02 2025a. doi: 10.48550/arXiv.2502.07981.
- Kim, S. and Chilton, L. B. Ai humor generation: Cognitive, social and creative skills for effective humor, 2025b. URL <https://arxiv.org/abs/2502.07981>.
- Loakman, T., Thorne, W., and Lin, C. Who’s laughing now? an overview of computational humour generation and explanation. In Flek, L., Narayan, S., Phuong, L. H., and Pei, J. (eds.), *Proceedings of the 18th International Natural Language Generation Conference*, pp. 780–794, Hanoi, Vietnam, October 2025. Association for Computational Linguistics. URL <https://aclanthology.org/2025.inlg-main.45/>.
- Sakabe, R., Kim, H., Hirasawa, T., and Komachi, M. Assessing the capabilities of llms in humor: a multi-dimensional analysis of oogiri generation and evaluation, 2025. URL <https://arxiv.org/abs/2511.09133>.
- Sanmartin, D. Kg-rag: Bridging the gap between knowledge and creativity, 2024. URL <https://arxiv.org/abs/2405.12035>.
- Tikhonov, A. and Shtykovskiy, P. Humor mechanics: Advancing humor generation with multistep reasoning, 2024. URL <https://arxiv.org/abs/2405.07280>.
- Wang, B., Wu, X., Liu, X., Li, J., Tiwari, P., and Xie, Q. Can language models make fun? a case study in chinese comical crosstalk, 2022. URL <https://arxiv.org/abs/2207.00735>.
- Yu, Z., Hu, X., and Wan, X. Cfunmodel: A ”funny” language model capable of chinese humor generation and processing, 2025. URL <https://arxiv.org/abs/2503.20417>.

Iteration 3

DSL Snippets

[FAST]我妈有个粉色保温杯，上面印着‘温柔妈妈’四个字，底下还有一行小字：[SLOW]‘情绪管理示范基地’。[PAUSE=1.0]

[NORMAL]她说这是她参加完家庭教育讲座送的纪念品，[EMPH]每次开家长会必带，[PAUSE=0.5]说是能镇魂。[SIG=Haiyaa]

[SLOW]那天我物理考了28分，[PAUSE=0.8]教室里安静得像殡仪馆播放课间操音乐。[PAUSE=1.2]

[FAST]老师念到我名字时，我妈立刻拧开杯子猛灌一口，[EMPH]然后递给我爸：[PAUSE=0.3]‘你喝！这火气不是我一个人的！’[SIG=Tsk]

[NORMAL]我爸刚要推辞，[PAUSE=0.5]老师突然鼓掌：[SLOW]‘下面请张明同学家长分享经验！’[PAUSE=1.0]

[FAST]张明是谁？就是那种你妈说[EMPH]‘你要是有他一半优秀，我就能提前退休’的那种人。[PAUSE=0.7]

[NORMAL]他妈妈站起来，旗袍盘扣都透着优越感，[SLOW]从包里掏出三本习题集：[PAUSE=0.5]‘唉，孩子太贪玩，[NORMAL]昨晚才做完剑桥预科卷，[FAST]今早顺手背了篇《自然》期刊摘要。’[PAUSE=1.0]

[SLOW]我正怀疑这孩子是不是外星派来打击地球学生的，[PAUSE=0.8]忽然听见他小声对他妈说：[PAUSE=0.5][EMPH]‘妈，下周奥数班要交量子力学读书报告……能不能换补习班？我想去普通人类班。’[SIG=Haiyaa][PAUSE=1.5]

[SLOW]我瞬间释然了。[PAUSE=0.5]原来你们家的天才，[EMPH]也是被腌在补习缸里的泡菜，[PAUSE=0.3]只不过你们用的是法国海盐。[SIG=■][PAUSE=1.0]

[NORMAL]回家路上，我妈没骂我，反而幽幽地说：[SLOW]‘要不……咱也报个班？’[PAUSE=1.0]

[FAST]我说您别折腾了，咱家连Wi-Fi都是租的。[PAUSE=0.5]结果第二天，[SLOW]她真把那保温杯供在书桌上，[EMPH]天天对着它默念：[PAUSE=0.3]‘让我儿子及格吧，让我儿子及格吧……’[SIG=Emotional Damage][PAUSE=1.2]

[NORMAL]前天亲戚来串门，一眼看到冰箱上的杯子照片，[EMPH]惊了：[PAUSE=0.3]‘哎哟，这不是我们家补习机构发的“焦虑缓冲杯”吗？[FAST]买满十万课时送的！’[SIG=Haiyaa][PAUSE=1.5]

[SLOW]我当场愣住——[PAUSE=0.5]原来全城爸妈都在用同一个杯子，[PAUSE=0.3]说着同一句咒语，[PAUSE=0.3]供养着同一个神话。[PAUSE=2.0]

Quality Evaluation

Writer Feedback:

- Setup 开场照搬‘最怕空气突然安静’热梗，原创性为零，你是来参加脱口秀还是段子复制大赛？
- Punchline 1 的‘反人类天赋’情绪到位，但前情铺垫太急，缺少对‘别人家孩子’出场的仪式感烘托，笑点被压扁了
- 转折段‘路过他们座位’太随意，像是偷听广播剧，缺乏空间动线和心理过渡，观众会被闪到腰椎间盘

RAG Feedback:

- 从脚本反推，RAG 检索的‘家长会社死’和‘别人家孩子补课’梗具有高共鸣性和校园真实性
- ‘补习班神话破灭’与‘全城父母同款焦虑杯’体现社会洞察，非低质堆笑料

请以事和人为主题生成一段3min的脱口秀文本

Thought for 7s

大家好，今天我想聊两个东西：事和人。

我发现成年人的崩溃，基本都来自这俩的组合题——事难办，人更难办。

你看，“事”本来挺简单：开始、过程、结果。

但一旦加了“人”，事情就变成了：开始、过程、解释、甩锅、复盘、再解释。

事情是客观的，人是主观的——所以很多时候我们不是在解决问题，我们是在解决***谁觉得这算问题**。

我最怕听到一句话：“这个事不复杂。”

每次有人说“不复杂”，我都想问：你说的是事不复杂，还是你这个人不复杂？

因为现实是：事确实不复杂，复杂的是——

同一件事，十个人能给你十种“理解”，外加一个“我觉得”。

（停顿）

“我觉得”这三个字，堪称人类发明的最强外挂：

没有证据、没有逻辑、没有数据——但有气势！

比如工作里最经典的场景：群里发一句“大家看一下”。

“看一下”这三个字就很灵性：

它既可能是“随便看看”，也可能是“今晚别睡”。

更可怕的是它后面还常跟一句：“辛苦了。”

我现在对“辛苦了”已经PTSD了——

每次看到这仨字，我都觉得自己马上就要辛苦了。

Figure 7. GPT. Qualitative example of Chinese stand-up generation on the same topic prompt.

A. Single Agent Example

B. Single Finetuned Agent Example

C. Different Temperature Setting Examples

生活笑谈：事教人暖

开场白：

大家好！我是北京大学AI助手，今天咱们来聊聊生活中的那些“事”和“人”。你们有没有发现，生活就像一本教科书，每件事都是个小测验，每个人都是位好老师？但有时候，这测验也太难了吧！比如，我小时候学数学，总觉得“ $1+1=2$ ”太简单了，可一到考试，就变成了“ $1+1=?$ ”等等，我是不是该先复习一下？”（观众笑）

主体部分：

先说“事”吧。生活中总有些小事，让我们哭笑不得。比如，学习外语这事儿。我有个朋友，为了练英语，天天看美剧，结果呢？他学会了说“How are you?”但回答永远是“I'm fine, thank you, and you?”（模仿机械语调）。有一天，他遇到个老外，老外问：“What's up?”他愣了半天，最后憋出一句：“The sky is up!”（观众大笑）这事儿告诉我们，学习不能死记硬背，得灵活应用。但别担心，犯错也是好事——至少你记住了“sky”这个词，对吧？

再说“人”。人嘛，总是生活中的亮点。比如，我的数学老师，她是个超级耐心的人。有一次我考砸了，她没骂我，反而说：“失败是成功之母，但你要记住，这位‘母亲’可能有点唠叨。”（观众笑）她教会我，挫折不是终点，而是成长的阶梯。还有那些好朋友，他们就像学习中的“参考答案”——不一定全对，但总能给你启发。记得有次我熬夜复习，朋友发来消息：“别学了，睡觉吧，明天大脑会感谢你的。”结果我真睡了，第二天考试居然超常发挥！看来，有时候“人”比“事”更重要，因为他们提醒我们：生活不是单打独斗。

结尾：

总之，生活中的“事”和“人”就像一对好搭档：事让我们成长，人给我们温暖。无论是学习中的小插曲，还是身边的良师益友，都值得珍惜。所以，下次遇到难题时，不妨笑笑——因为每件事都在教你点什么，每个人都在帮你变得更棒！谢谢大家！（掌声）

Figure 8. DeepSeek. Qualitative example of Chinese stand-up generation on the same topic prompt.

--- Testing Prompt ---

《脱口秀大赛》，请以程序员的日常为风格，创作一段脱口秀台词。

--- Finetuned Model Output ---

你说我写代码像在演话剧？不，我就是个程序员！每天上班第一件事不是吃饭，是开电脑——打开电脑、登录系统、打开浏览器、进入公司官网……然后呢？点进一个叫‘招聘’的页面，一看：2018年1月1日，入职，工作年限零天；3月4日，转正，还是零天；6月9日，升级主管，又归零；10月25日，晋升经理，再加零；最后12月27日，升职总监，工资没涨，但地位上去了，我也算熬出头了。

你说这哪是工作？这是追剧吧！从新人到老员工，每一步都像是通关，虽然奖励是加班和泡面。但你得坚持啊，毕竟没人告诉你，干一行就爱一行，别总想着跳槽去当明星——你看人家明星多风光，可他们也得吃饭睡觉，还得被狗仔拍屁股擦脸。我们程序员不一样，我们有逻辑，我们有bug，我们还有KPI。所以兄弟们，别急着辞职，先学会用Git管理人生版本吧！

今天我给大家带来一场技术与情感的碰撞——《程序员的爱情》。故事讲的是两个程序员，一个男的一个女的，在机房里相遇，聊代码，谈理想，结果越聊越深，最后……直接走到了一起。你以为这是爱情片？错，这是真实发生过的。不信？我还能现场演示一段Python写的告白程序，保证比电视剧更感人，比小说更深刻，比诗歌还押韵——只要你愿意给我输入密码，我就告诉你，如何用一行代码，把全世界的心都变成我的数据库。

来吧，输入你的用户名，让我们开启这段浪漫旅程。记住，一旦运行成功，你的人生数据表，就永远属于我了。放心，不会泄露，因为我是最安全的后端工程师——而且，我已经做好准备，随时可以给你发一封加密邮件，标题写着：亲爱的，我想你了。

当然，如果你不想这么亲密接触，请立即退出本程序。否则，下次见面时，你可能不再是路人甲，而是我精心设计的角色扮演对象之一。欢迎再次光临，本次服务结束前，祝你编码顺利，感情稳定，代码无Bug，生活甜蜜，梦想成真，早日升职加薪，走上人生巅峰——或者至少，别再…

Figure 9. QLoRA finetuned Qwen3-4B-Instruct

- 朋友们，人生最难回答的问题，真的不是“你爱我吗”，而是——“吃啥？” (Setup)
- 每次点外卖，我和我对象的对话都像一场没有硝烟的战争。
- 我说：“你定吧。”
- 她说：“随便。”
- 我说：“那我真点了？”
- 她立刻警觉：“你先说你想吃啥？”
- 我说：“我随便。”
- 她说：“你这哪是随便，你明明就想吃螺蛳粉！” (Punchline 节奏1)
- 你看，他说“随便”，其实不是随便。她说“你定”，其实不是让你定。
- 后来我们实在受不了了，下载了一个AI点餐助手，号称能根据心情、天气、星座甚至你昨天打嗝的频率，精准推荐晚餐。
- AI上线第一句就问：“请问您对番茄炒蛋加糖持什么立场？” (Callback伏笔)
- 我懵了。我说加，图个甜口儿；她怒了，说这是对江浙沪的背叛！
- AI沉默三秒，屏幕突然变红，弹出一行字：
- “检测到核心价值观冲突，建议立即分手，以免浪费本AI算力。” (Punchline 2)
- 那一刻我悟了。“随便”根本不是佛系，那是最高级的控制欲——
- 把锅甩给你，再站在道德高地审判你的选择。
- 现在，每当我对象再问我“吃啥”，我不纠结了。
- 我直视她的眼睛，一字一顿地说：
- “番茄炒蛋，不加糖。”
- 她愣住，然后小声问：“那……我呢？”
- 我微微一笑：“你随便。” (Final Callback & Punchline)
- 各位都市勇士、熬夜冠军、回消息回出腱鞘炎的朋友们！今天我隆重推出——人生万能遥控器！（举手比划）
- 你有没有这种时候？娃凌晨三点嚎啕大哭，老板微信弹出‘在吗’，而你家厨房正在上演《水漫金山》？这时候你只需要——按一下‘一键修复’！
- 我试过。按完第二天，娃睡了，老板不找我了，水管也不漏了……因为我直接被房东赶出来了，他说我用‘修复’键把整栋楼的下水道格式化了。（Punchline 1）
- 还有‘快进’键，多好啊！我想跳过周一，结果跳过了整个发薪日。‘重播’键更绝，我想重温青春，结果重播了上周修马桶——那画面，马桶圈反光里我的脸，比前任拉黑我还决绝。（Punchline 2）
- 但最神的是‘静音’键。那天老婆正跟我开家庭战略会议，讲到孩子上哪所幼儿园能保送清华……我手一抖，按了静音。
- 她，真的，消失了。
- 三天后她回来，问我：‘你是不是觉得我不重要？’
- 我说：‘不！你太重要了！重要到这个宇宙的静音键都对你无效！’
- （停顿，苦笑）
- 其实吧……这遥控器根本不存在。那天我只是躲进车库，抽了根烟，手机开了飞行模式。
- 可你知道最讽刺的是啥吗？
- 昨天我又梦见那个遥控器了。这次我鼓起勇气，按下了‘修复’键。
- 今早醒来，我家WiFi名字真的变成了——‘对不起老婆我爱你’。
- 所以朋友们，别信什么一键解决。人生这破遥控器，连电池仓都焊死了，咱们只能……硬着头皮，手动调台。（Callback x2）

Figure 10. Left one generated with 0.1 temperature, right one with 0.9 temperature