

# Critical Analysis: Synchronous Faithfulness Monitoring for Trustworthy Retrieval-Augmented Generation

Kai-Yu Lu

## 1 Methodological Strengths

### 1.1 Comprehensive Benchmark Design

A central strength of the paper lies in the construction of a multi task benchmark that covers four long form retrieval augmented generation (RAG) settings: question answering, summarization, data to text generation, and biography generation. The benchmark spans six datasets: three from RAGTruth (QA, Summ, Data2txt), the factscore benchmark (FS), and two newly constructed biography datasets (F-100 and F-100-anti). This design has several advantages.

First, the benchmark covers both classic RAG settings such as QA and summarization, and more challenging entity centric biography generation, where parametric knowledge is strong and can conflict with retrieved context. F-100 uses aligned Wikipedia contexts, while F-100-anti intentionally misaligns context via entity substitution, thereby stressing the model’s ability to resist generated content that ignores or contradicts the provided evidence.

Second, the paper standardizes sentence level faithfulness labels across datasets. For RAGTruth tasks, annotated baseless and conflict spans are used to mark unfaithful sentences. For FS, F-100, and F-100-anti, an automatic pipeline decomposes outputs into propositions and uses an AutoAIS based fact checker with lexical mapping back to sentences. This enables consistent evaluation of sentence level detectors across diverse tasks and base models.

Third, the authors report detailed dataset statistics, such as average numbers of context sentences, output sentences, and proportions of faithful sentences and faithful instances for Llama 2 7B Chat, with analogous patterns for Mistral 7B Instruct. This transparency clarifies how challenging each task is for faithfulness tracking and grounds subsequent AUROC comparisons.

### 1.2 Feature Based Synchronous Monitoring

The proposed monitor, SYNCHECK, is methodologically appealing because it integrates four complementary families of decoding time features that correspond to distinct failure modes of RALMs:

- **Likelihood** features capture knowledge gaps through minimum and length normalized sequence likelihood, aligning with the intuition that low probability outputs often indicate unsupported content.
- **Uncertainty** features combine token level entropy and local intrinsic dimension of intermediate activations. This goes beyond raw probabilities and attempts to characterize unstable or high dimensional representation manifolds associated with unfaithful behavior.

- **Context influence** features compare token distributions conditioned on inputs with and without retrieved context via Kullback–Leibler divergence, directly targeting over reliance on parametric knowledge and under use of retrieval.
- **Semantic alignment** features incorporate an external entailment style checker to measure whether generated sentences are semantically supported by the context.

This design explicitly reflects the authors’ conceptual hierarchy of untrustworthy behavior: insufficient joint knowledge, ignoring context, and misusing context. By grounding the feature set in these categories, the detector is better aligned with the target property of context faithfulness than methods that rely solely on likelihood or lexical overlap.

### 1.3 Lightweight Aggregator and Strong Baseline Comparisons

Another methodological strength is the choice of simple aggregators for the feature set, namely logistic regression, XGBoost, and a small multilayer perceptron. The paper shows that SYNCHECK with logistic regression already achieves high AUROC, only slightly below the MLP variant. This implies that the main gains come from feature design rather than complex classifier architectures, enhancing interpretability and deployability.

The evaluation compares SYNCHECK against a wide range of baselines:

- SPANEXTRACT with GPT-4-Turbo and fine tuned Llama 2 13B as span level hallucination detectors.
- Lexical alignment models such as AlignScore and MiniCheck.
- System quality control signals such as CRITICTOK from Self-RAG and the FLARE likelihood threshold.

Performance is reported as AUROC for each dataset and model, with averages across tasks. SYNCHECK<sub>MLP</sub> achieves average AUROC around 0.831 for Llama 2 7B Chat and 0.867 for Mistral 7B Instruct, improving prior methods by 4 to 35 points depending on the task. The inclusion of strong external baselines and detailed per task results strengthens the claims regarding improved faithfulness detection.

### 1.4 Faithfulness Oriented Decoding Protocol

The Faithfulness Oriented Decoding (FOD) algorithm is a methodologically coherent second component that builds directly on SYNCHECK. It uses a two stage protocol:

1. **Greedy search with backtracking:** generate sentences greedily until the first sentence whose SYNCHECK score falls below a threshold  $\tau_1$ , then backtrack and keep the prefix up to the last faithful sentence.
2. **Faithfulness guided beam search:** starting from the faithful prefix, run beam search where candidate continuations are sampled and immediately pruned if their SYNCHECK scores fall below a threshold  $\tau_2$ . The search keeps the top  $K$  beams according to aggregated faithfulness.

This protocol transforms sentence level monitoring into a concrete decoding constraint, providing a mechanism to bound minimum sentence level faithfulness. The authors evaluate FOD on the same six datasets using two response level metrics: proposition based faithfulness and informativeness (number of propositions). They compare FOD with greedy decoding, abstention, reranking, and context aware decoding (CAD), showing consistent gains in faithfulness and favorable trade offs between faithfulness and informativeness.

## 1.5 Systematic Ablations and Generalization Studies

The methodology also includes several systematic analyses:

- **Feature ablation:** removing each feature family in turn shows that all four contribute to AUROC, with semantic alignment most critical and context influence second.
- **Cross task transfer:** training SYNCHECK on one task and evaluating on another demonstrates that detectors trained on FS or Data2txt transfer reasonably well to QA and other tasks.
- **Cross model transfer:** training SYNCHECK on Llama 2 7B Chat and evaluating on other Llama Chat models or Mistral 7B Instruct still yields strong AUROC, suggesting that features remain informative under model shift.

These studies provide evidence that the proposed features generalize across tasks and base models without heavy re training, which is essential for practical trustworthiness tools.

## 2 Key Limitations

### 2.1 Scope and Representativeness of the Benchmark

Although the benchmark spans six datasets, its coverage of task types and domains remains limited in several respects.

First, all tasks are text only and fall into four categories: QA, news summarization, Yelp style data to text, and Wikipedia style biographies. Critical real world RAG applications such as multi turn conversational agents, code generation, legal or financial advice, and multi modal retrieval are not included. As a result, the generalization of SYNCHECK and FOD to these settings is not empirically substantiated.

Second, the biography datasets F-100 and F-100-anti focus on highly popular entities, where parametric knowledge is strong. While this is a natural stress test for context conflicting hallucinations, it underrepresents less popular entities and domains where retrieved knowledge is sparse or noisy, or where contexts are long, heterogeneous, and partially irrelevant.

Third, the benchmark is monolingual and restricted to English, which leaves cross lingual and multilingual faithfulness entirely unexplored.

### 2.2 Quality and Noise in Automatic Labeling

For FS, F-100, and F-100-anti, sentence level labels are derived from an automatic pipeline that combines propositionization, automatic fact checking with AutoAIS, and lexical alignment back to sentences. Each stage introduces potential errors:

- Propositionization may split or group factual content in ways that misrepresent the underlying sentence semantics.
- Automatic fact checking can misjudge borderline or context dependent claims, especially when nuanced temporal or contextual qualifiers are present.
- Lexical mapping may fail when sentences paraphrase propositions in non trivial ways.

The paper does not provide a detailed analysis of label noise induced by this pipeline or manual audits of sentence level annotations. Since these labels are used both for training SYNCHECK and for evaluation, accumulated noise may bias AUROC estimates and obscure failure modes.

## 2.3 Granularity of Faithfulness Assessment

The entire framework treats one sentence as one segment for monitoring and decoding. This design is practical but introduces several limitations:

- Many hallucinations are clause level or span level rather than sentence level. Labeling an entire sentence as unfaithful when only a clause is problematic discards the faithful part of the sentence in evaluation and training.
- Conversely, a sentence that contains an unfaithful clause but is dominated by faithful content can be wrongly judged as faithful if the fact checking pipeline fails to identify the problematic proposition.
- FOD uses sentence level thresholds to backtrack and prune beams, which may be too coarse in scenarios where segment granularity needs to align with finer semantic units.

The paper acknowledges segment granularity as a limitation but does not explore alternative units such as clauses or multi sentence segments.

## 2.4 Evaluation of Intervention Side Effects

The faithfulness intervention experiments focus on two metrics: faithfulness and informativeness, both defined at proposition level. Although these metrics are appropriate for factual consistency, they do not capture several important aspects:

- Coherence and discourse structure of generated texts under FOD are not evaluated. Early backtracking and beam search may introduce repetition or awkward transitions.
- Human preferences and task specific utility are not measured. An answer that is slightly less faithful but significantly more helpful to users is not distinguished from an answer that is faithful but unhelpful.
- The paper does not report qualitative user studies or detailed error case analyses of FOD outputs.

Without such evaluations, the broader impact of FOD on response quality remains only partially understood.

## 2.5 Computational Cost and Deployment Considerations

SYNCHECK requires several additional computations beyond standard decoding:

- Computing context influence involves an extra forward pass per token without context.
- Semantic alignment requires running an external entailment style model for each generated sentence.
- FOD performs sampling and beam search with faithfulness checks for each candidate continuation.

The limitations section briefly discusses latency and notes that FOD with sample size one has similar cost to CAD, and that parallelization across multiple GPUs can mitigate overhead. However, the main text does not present concrete wall clock measurements, breakdown of overhead per component, or analysis of cost as a function of beam size and sample size. The absence of such data makes it difficult to assess practical feasibility in resource constrained or latency sensitive deployment scenarios.

## 3 Technical Bottlenecks

### 3.1 Dependence on Internal Model Signals

SYNCHECK relies heavily on internal model signals such as token probabilities, token level entropy, and intermediate layer activations. This imposes several constraints.

First, access to token level log probabilities and hidden states is required to compute likelihood, uncertainty, and local intrinsic dimension. This is feasible for open source models like Llama 2 and Mistral, but less straightforward for closed models that expose only sampled tokens or limited log probability information.

Second, context influence requires computing predictive distributions with and without retrieved context for the same prefix. This doubles the number of forward passes and presupposes that the model API permits full control over input context, which may not hold for some production systems that fuse retrieval and prompting internally.

Surrogate model experiments show that SYNCHECK trained on one model can judge outputs from another model; however, this still assumes the availability of a surrogate with similar architecture and full access to its internals.

### 3.2 Information Bottlenecks in Aggregation

The feature aggregator maps a relatively rich set of decoding time statistics into a single scalar faithfulness score per sentence. This introduces at least two information bottlenecks:

- Distinct failure modes that share similar scalar scores are no longer distinguishable during decoding. For instance, sentences with high entropy and low context influence can be conflated with sentences that have moderate entropy and moderate semantic misalignment.
- FOD uses simple averaging of sentence level scores over prefixes and beams to rank candidate outputs. This aggregation does not distinguish scenarios where a single highly unfaithful sentence is embedded in an otherwise faithful answer from scenarios where faithfulness is uniformly moderate.

These bottlenecks limit the expressiveness of the monitoring signal and constrain the kinds of interventions FOD can apply.

### 3.3 Decoding Trade Offs and Hyperparameter Sensitivity

FOD introduces several hyperparameters, including backtrack threshold  $\tau_1$ , pruning threshold  $\tau_2$ , beam size  $K$ , and per step sample size  $S$ . The main experiments fix  $\tau_1 = 0.7$ ,  $\tau_2 = 0.85$ ,  $K = 2$ , and  $S = 6$  across all tasks and models. While this demonstrates robustness to some extent, it also reveals a bottleneck:

- The method implicitly assumes that a single operating point balances faithfulness and informativeness across diverse domains and base models, which may not hold in higher risk or more heterogeneous settings.
- Fine grained tuning of thresholds for different user risk profiles or applications is not investigated in the main text. Appendix analyses of FOD hyperparameters are limited and do not explore dynamic or adaptive thresholding strategies.

This constrains the flexibility of FOD and may hinder its adoption in applications that require precise control over abstention rates or faithfulness guarantees.

### 3.4 Limited Integration with Retrieval and Knowledge Sources

Although SYNCHECK monitors context influence and detects over dominance of parametric knowledge, the intervention protocol does not modify retrieval itself. Retrieval quality and coverage are taken as fixed:

- When faithfulness scores degrade due to poor or irrelevant context, FOD can only prune or backtrack on generation; it cannot trigger re retrieval, context expansion, or context filtering.
- This contrasts with methods such as FLARE or iterative retrieval that integrate retrieval decisions into decoding. In the current framework, retrieval and decoding remain loosely coupled.

This separation creates a bottleneck: many failure cases may require adaptation of both retrieval and decoding, but only the latter is directly controlled.

## 4 Research Implications

### 4.1 Understanding the Limits of Existing Faithfulness Signals

The benchmark results reveal that widely used signals for quality control in RAG pipelines have limited effectiveness for sentence level faithfulness detection:

- CRITICTOK and FLARE, which rely on critique tokens and least confident token likelihood, achieve AUROC values around 0.6 across several tasks.
- External lexical alignment models such as AlignScore perform well on some tasks but generalize poorly to data to text tasks that require fine grained numerical consistency.

These findings indicate that existing ad hoc signals are insufficient for robust faithfulness monitoring in long form RAG. For the broader field, this suggests that faithfulness is a distinct property that requires dedicated modeling rather than being treated as a byproduct of accuracy or calibration.

### 4.2 Evidence for Rich Meta Information in Decoding Dynamics

The success of SYNCHECK suggests that decoding dynamics of RALMs contain rich meta information about the trustworthiness of outputs:

- Local intrinsic dimension and entropy capture forms of epistemic uncertainty that correlate with unfaithful behavior.
- Divergences between distributions with and without context quantify how strongly retrieval influences predictions, providing a mechanistic view of context usage.
- Combining these with semantic alignment yields a detector that generalizes across tasks and models.

This strengthens the view that internal states of LLMs can provide reliable self diagnostic signals, which is relevant for broader work on confidence estimation, selective generation, and self evaluation.

### **4.3 Decoding as a Vehicle for Trustworthiness Control**

FOD demonstrates that decoding can be explicitly structured around trustworthiness objectives rather than solely optimizing likelihood or task accuracy. The algorithm provides:

- A constructive example of how sentence level monitors can be integrated into search to enforce minimum faithfulness thresholds.
- An empirical demonstration that such interventions can improve both faithfulness and informativeness compared to common baselines such as abstention and CAD.

This has implications for other domains where trustworthiness constraints are important, such as safety critical question answering, legal drafting, or medical summarization, where decoding algorithms may need to embed structural guarantees.

### **4.4 Modular Trustworthiness Infrastructure**

Cross task and cross model generalization of SYNCHECK indicates that trustworthiness monitors can be designed as modular components that are not tightly coupled to a single task or base model. A detector trained on one RAG setting can often be reused in others with acceptable degradation, and a surrogate model can support monitoring of black box models.

This modularity is promising for building layered trust infrastructures where detection and decoding components can be swapped or upgraded independently from base models and retrievers.

## **5 Potential Research Directions**

### **5.1 Richer and More Diverse Benchmarks**

Future work can extend the benchmark along several dimensions:

- Include multi turn conversational RAG tasks where context faithfulness must be maintained over dialogue histories and not just single questions.
- Incorporate multi modal RAG tasks that involve images, tables, or code, assessing whether decoding time signals still correlate with faithfulness when modalities interact.
- Add domains with high stakes such as medicine, law, or finance, where error types and tolerance levels differ significantly from generic QA or biography generation.
- Construct multilingual and cross lingual RAG benchmarks to evaluate whether SYNCHECK style features transfer across languages.

Such expansions would clarify the limits of current monitoring and decoding methods and reveal domain specific challenges.

### **5.2 Improved Labeling Pipelines and Ground Truth Audits**

Improving the quality of faithfulness labels is an important direction:

- Design propositionization models specifically tuned for RAG outputs and domain specific discourse structures.

- Develop hybrid pipelines that combine automatic fact checking with targeted human verification on difficult or ambiguous cases, thereby benchmarking detector performance under lower label noise.
- Create error taxonomies for unfaithfulness, distinguishing baseless hallucinations, context conflicts, partial misinterpretations, and subtle evidential overstatements, and annotate data accordingly.

These steps would support more nuanced evaluations and more targeted training of detectors.

### **5.3 Finer Grained Monitoring and Structured Interventions**

Moving beyond sentence level monitoring offers another avenue:

- Extend SYNCHECK to operate at clause or span level, possibly aligned with proposition boundaries, to localize unfaithful content more precisely.
- Couple decoding with local editing operations, such as deleting or rewriting only unfaithful spans, while retaining faithful parts of sentences.
- Integrate discourse level constraints to ensure that corrections do not disrupt coherence or introduce contradictions across sentences.

Such structured interventions would refine the balance between faithfulness and readability.

### **5.4 Joint Optimization of Retrieval and Decoding**

A natural extension is tighter integration between retrieval and decoding:

- Use low context influence scores as triggers for re retrieval, context expansion, or context filtering, thereby closing the loop between knowledge acquisition and generation.
- Learn retrieval policies that explicitly optimize downstream faithfulness scores, not just relevance or similarity.
- Explore multi stage pipelines where retrieval and decoding are jointly optimized through reinforcement or bandit style feedback using faithfulness metrics.

This would address failure modes where unfaithfulness originates more from retrieval deficiencies than from decoding behavior.

### **5.5 Adaptive Thresholding and Risk Sensitive Decoding**

Current FOD uses fixed thresholds across tasks and models. Future work can:

- Develop calibration procedures that map SYNCHECK scores to explicit estimates of error probabilities, enabling thresholds tied to desired risk levels.
- Implement dynamic thresholding schemes that adjust  $\tau_1$  and  $\tau_2$  based on user specified risk preferences, domain, or query difficulty.
- Study policies that mix abstention, FOD style intervention, and fallback to human experts or secondary systems when risk is high.

These strategies would enable more flexible deployment in contexts with varying tolerance for unfaithfulness.

## 5.6 Learning Richer Meta Models

Finally, research can explore more expressive meta models for faithfulness:

- Replace or augment hand crafted features with learned representations from dedicated meta models that ingest full decoding traces, attention patterns, or gradients.
- Train joint models that predict multiple attributes, such as factuality, robustness to adversarial contexts, and sensitivity to context perturbations, enabling multi objective decoding.
- Investigate self supervised training signals derived from controlled perturbations of context, such as synthetic contradictions or context removal, to scale faithfulness monitoring without heavy annotation.

These approaches could further leverage the rich internal structure of RALMs.

## 6 Conclusion

The paper makes a substantial contribution to the study of trustworthiness in retrieval augmented generation by introducing SYNCHECK, a feature based synchronous faithfulness monitor, and FOD, a decoding algorithm that leverages this monitor to improve faithfulness while retaining informativeness. Methodologically, the work is distinguished by a multi task benchmark, a principled feature design anchored in decoding dynamics, and systematic comparisons with strong baselines.

At the same time, the study exhibits limitations in benchmark scope, reliance on automatic labeling, coarse sentence level granularity, incomplete evaluation of intervention side effects, and limited analysis of computational overhead. These limitations point directly to technical bottlenecks in monitoring and decoding for trustworthy RAG.

The most promising research directions include expanding and refining benchmarks, improving labeling quality, developing finer grained and more structured interventions, jointly optimizing retrieval and decoding, designing adaptive risk sensitive decoding protocols, and learning richer meta models that exploit full decoding traces. Progress along these directions would further advance the goal of building RAG systems that are not only powerful but also reliably faithful to the evidence they are given.