

Of course. Here is a detailed breakdown of the HellaSwag document's core logic, mathematical models, and symbolic representations, followed by a new, generic framework for next-event prediction systems.

(1) Pure Logic of the Document

The document's core argument is structured as follows:

* **Problem:** Despite models like BERT achieving near-human performance on the SWAG dataset, this success is misleading. Models are not learning genuine commonsense reasoning but are instead "rapid surface learners" that exploit subtle dataset-specific biases and stylistic patterns.

* **Hypothesis:** A more robust benchmark can be created using **Adversarial Filtering (AF)**, a data collection paradigm that uses state-of-the-art models to iteratively select machine-generated wrong answers that are maximally confusing for those same models, while remaining nonsensical to humans.

* **Method:**

1. **Source Diverse Contexts:** Use contexts from two domains: ActivityNet (video captions) and WikiHow (how-to articles). The latter provides longer, more complex contexts.

2. **Generate Negative Endings:** Use a powerful language model (GPT) to generate a massive number of potential incorrect endings for each context.

3. **Adversarial Filtering:** Iteratively train a discriminator (BERT-Large) on a subset of the data and use it to identify and replace "easy" negative endings in the held-out set with more challenging, "adversarial" ones from the generator pool. This process converges when the discriminator can no longer reliably identify the correct ending.

4. **Human Validation:** Use human annotators to filter out any machine-generated endings that are accidentally plausible, ensuring the final dataset is trivial for humans.

* **Evaluation:**

- * **Primary:** Show that the resulting dataset, HellaSwag, is extremely challenging for all state-of-the-art models (<48% accuracy) while being easy for humans (>95% accuracy).

- * **Secondary:** Demonstrate that this difficulty persists even on in-domain examples and worsens on "zero-shot" categories, proving that models are not learning generalizable commonsense.

- * **Conclusion:** HellaSwag serves as a genuinely challenging testbed for commonsense reasoning. The paper advocates for a future where benchmarks **co-evolve** with model capabilities through adversarial methods, ensuring they present ever-harder challenges and accurately measure progress on the underlying **task** rather than mere dataset fitting.

Revised Mathematical Framework for HellaSwag

1. Formal Definition of the Adversarial Filtering Process

Let us define the adversarial filtering as a **minimax optimization problem** over probability spaces:

Definition 1 (Adversarial Filtering Game):

Let:

- \mathcal{C} be the space of contexts

- \mathcal{E} be the space of endings

- \mathcal{D} be the distribution over (context, correct ending) pairs: $(c, e^+) \sim \mathcal{D}$

- $G: \mathcal{C} \rightarrow \mathcal{P}(\mathcal{E})$ be the generator function producing candidate negative endings

- \mathcal{F} be the class of discriminator functions $f: \mathcal{C} \times \mathcal{E} \rightarrow \mathbb{R}$

The adversarial filtering process solves:

$$\min_{\mathcal{D}_{\text{bench}}} \max_{\{f \in \mathcal{F}\}} [\mathbb{E}_{\{(c, e^+) \sim \mathcal{D}_{\text{bench}}\}} [\log f(c, e^+)] - \mathbb{E}_{\{(c, e^+) \sim \mathcal{D}_{\text{bench}}, e^- \sim G(c)\}} [\log f(c, e^-)]]$$

where $\mathcal{D}_{\text{bench}}$ is the benchmark distribution being constructed.

2. Iterative Process as Fixed-Point Iteration

Theorem 1 (AF as Contraction Mapping):

Define the adversarial update operator $T: \mathcal{P}(\mathcal{C} \times \mathcal{E}) \rightarrow \mathcal{P}(\mathcal{C} \times \mathcal{E})$ as:

$$T(\mathcal{D}) = \{(c, e^+) \in \mathcal{D} : f_D(c, e^+) < \tau\} \cup \{(c, e^+) : e^+ \text{ replaced with more adversarial } e^-\}$$

where f_D is the optimal discriminator trained on \mathcal{D} and τ is a threshold.

If T is a contraction mapping on the metric space $(\mathcal{P}(\mathcal{C} \times \mathcal{E}), d)$ for some metric d , then the AF process converges to a unique fixed point.

Proof Sketch:

We need to show:

1. **Completeness:** $(\mathcal{P}(\mathcal{C} \times \mathcal{E}), d)$ is a complete metric space

2. **Contraction:** $d(T(\mathcal{D}_1), T(\mathcal{D}_2)) \leq k \cdot d(\mathcal{D}_1, \mathcal{D}_2)$ for some $k < 1$

For the Wasserstein metric W_1 , under Lipschitz conditions on f and G , we can establish contraction.

3. Convergence Analysis

Proposition 1 (Rate of Convergence):

Under the contraction mapping framework, the AF process exhibits linear convergence:

$$d(\mathcal{D}_{\{n+1\}}, \mathcal{D}^*) \leq k^n \cdot d(\mathcal{D}_1, \mathcal{D}^*)$$

where \mathcal{D}^* is the fixed-point distribution and k is the contraction coefficient.

Definition 2 (ϵ -Adversarial Equilibrium):

A benchmark \mathcal{D} is in ϵ -adversarial equilibrium if:

$$|\mathbb{E}_{\{(c, e^+) \sim \mathcal{D}\}} [f(c, e^+)] - \mathbb{E}_{\{(c, e^+) \sim \mathcal{D}, e^- \sim G(c)\}} [f(c, e^-)]| < \epsilon \quad \forall f \in \mathcal{F}$$

4. Information-Theoretic Formulation

Definition 3 (Adversarial Information Gap):

The quality of the benchmark can be measured by:

$$\Delta_I(\mathcal{D}) = I(C; E^+) - \max_{\{f \in \mathcal{F}\}} I_f(C; E^-)$$

where I is mutual information and I_f is the f -weighted mutual information.

****Theorem 2** (Optimality Condition):**

The AF process minimizes the adversarial information gap $\Delta_I(\mathcal{D})$ while maintaining $I(C; E^+)$ near the human level.

5. Stochastic Process Formulation

Model the AF iterations as a Markov chain $\{\mathcal{D}_t\}_{t \geq 0}$ with state space $\mathcal{P}(\mathcal{C} \times \mathcal{E})$.

****Definition 4** (AF Markov Chain):**

The transition kernel P is defined by:

$$P(\mathcal{D}_{t+1} = \mathcal{D}' \mid \mathcal{D}_t = \mathcal{D}) = \mathbb{P}(T(\mathcal{D}, \xi) = \mathcal{D}')$$

``

where ξ represents the randomness in discriminator training and negative sampling.

****Proposition 2** (Ergodicity):**

Under mild conditions on G and F , the AF Markov chain is ergodic and converges to a stationary distribution.

6. Geometric Characterization of the "Goldilocks Zone"

****Definition 5** (Adversarial Manifold):**

Let $M \subset \mathcal{P}(\mathcal{C} \times \mathcal{E})$ be the manifold of distributions where:

``

$$\mathbb{E}_{\{(c, e^+) \sim \mathcal{D}\}}[f(c, e^+)] \approx 1/N \quad \forall f \in F$$

``

but human accuracy remains high.

****Theorem 3** (Existence of Goldilocks Zone):**

If $\dim(F) < \infty$ and G is sufficiently expressive, then M is a non-empty submanifold of codimension $\dim(F)$.

7. Complete Mathematical Reframing

(2) Core Mathematical Models - Revised

1. Formal Adversarial Filtering Framework

Let $(\Omega, \mathcal{F}, \mathbb{P})$ be a probability space supporting:

- Context distribution: $C \sim \mu_{\mathcal{C}}$
- Correct ending map: $E^+: \mathcal{C} \rightarrow \mathcal{E}$
- Generator process: $G: \mathcal{C} \rightarrow \mathcal{P}(\mathcal{E})$
- Discriminator class: $F \subset \{f: \mathcal{C} \times \mathcal{E} \rightarrow \mathbb{R}\}$

****Definition 6** (Adversarial Value):**

``

$$V(\mathcal{D}, f) = \mathbb{E}_{\{(c, e^+) \sim \mathcal{D}\}}[\log \sigma(f(c, e^+))] + \mathbb{E}_{\{(c, e^+) \sim \mathcal{D}, e^- \sim G(c)\}}[\log(1 - \sigma(f(c, e^-)))]$$

``

where σ is the sigmoid function.

The AF process solves:

$$\min_{\mathcal{D} \in \mathfrak{D}} \max_{f \in F} V(\mathcal{D}, f)$$

where \mathfrak{D} is the set of feasible benchmark distributions.

2. Convergence Theorem

Theorem 4 (AF Convergence):

Assume:

1. F is compact and convex
2. G is Lipschitz continuous
3. The loss is strongly convex in f

Then the alternating updates:

$$f_{t+1} = \operatorname{argmax}_{f \in F} V(\mathcal{D}_t, f)$$

$$\mathcal{D}_{t+1} = \operatorname{argmin}_{\mathcal{D} \in \mathfrak{D}} V(\mathcal{D}, f_{t+1})$$

converge to a Nash equilibrium.

Proof: Apply Sion's minimax theorem and Banach fixed-point theorem.

3. Complexity Analysis

Definition 7 (Adversarial Complexity):

$$CAF(\epsilon) = \min\{T : |V(\mathcal{D}_T, f_T) - V^*| < \epsilon\}$$

where V^* is the game value.

Proposition 3 (Complexity Bound):

Under Lipschitz conditions, $CAF(\epsilon) = O(\log(1/\epsilon))$.

8. Practical Algorithm with Mathematical Guarantees

Algorithm 1 (Provably Convergent AF):

Input: Initial dataset \mathcal{D}_0 , generator G , tolerance ϵ

Output: Benchmark distribution \mathcal{D}^*

1. for $t = 0, 1, 2, \dots$ do
2. Train $f_t = \operatorname{argmax}_{f \in F} V(\mathcal{D}_t, f)$
3. Compute gradient: $\nabla_{\mathcal{D}} V(\mathcal{D}_t, f_t)$
4. Update: $\mathcal{D}_{t+1} = \Pi_{\mathfrak{D}}(\mathcal{D}_t - \eta_t \nabla_{\mathcal{D}} V(\mathcal{D}_t, f_t))$
5. if $\|\mathcal{D}_{t+1} - \mathcal{D}_t\| < \epsilon$ then break
6. end for
7. return \mathcal{D}_{t+1}

Theorem 5 (Convergence of Algorithm 1):

If step sizes η_t satisfy Robbins-Monro conditions, then Algorithm 1 converges almost surely to a local minimax equilibrium.

9. Statistical Testing Framework

Definition 8 (Adversarial Quality Metric):

$$Q(\mathcal{D}) = \inf_{\{f \in F\}} [ACC_human(\mathcal{D}) - ACC_f(\mathcal{D})]$$

where ACC measures accuracy.

Proposition 4 (Hypothesis Test):

We can test $H_0: Q(\mathcal{D}) \leq \delta$ against $H_1: Q(\mathcal{D}) > \delta$ using the empirical process:

$$\sqrt{n}(Q_n(\mathcal{D}) - Q(\mathcal{D})) \rightarrow \mathbb{G}$$

where \mathbb{G} is a Gaussian process.

(3) Symbolic Logic

Review of the Symbolic Logic Section

The symbolic logic in the HellaSwag_analysis.pdf is presented in Section (3) on page 5. It provides a basic formalization of the task as multiple-choice QA over situational contexts, including dataset structure, task function, and the adversarial filtering (AF) goal. This is concise and grounded in the original HellaSwag paper, emphasizing the contrast between model and human performance.

Strengths:

- Clearly defines the core elements (context $\langle c \rangle$, endings $\langle E \rangle$, correct ending $\langle e^+ \rangle$) using set notation.
- Ties directly to the AF process by formalizing the goal as minimizing model accuracy while maintaining high human accuracy, which aligns with the document's hypothesis.
- Practical and accessible, avoiding over-abstraction.

Weaknesses:

- **Lack of Depth in Reasoning Modeling**: It doesn't symbolize the commonsense aspects (e.g., physical, social, or temporal reasoning) that make HellaSwag challenging. The original paper implies categories like script knowledge or causality, but these aren't formalized (e.g., no predicates or entailment relations).
- **No Handling of Negatives**: The AF goal is stated probabilistically, but there's no symbolic way to represent why negatives are "adversarial" (e.g., via logical violations or plausibility measures).
- **Missing Integration**: It doesn't connect to the mathematical framework (e.g., minimax game in Definition 1 or Goldilocks Zone) or the generic framework (e.g., Generator/Discriminator roles). No uncertainty, causality, or lattice structures as in related frameworks.
- **No Complexity or Provability**: Lacks theorems on computational complexity (e.g., deciding plausibility) or algorithms for symbolic evaluation, making it less rigorous than the mathematical sections.
- **Asymmetry**: Focuses on the task but not on symbolic validation (e.g., how to prove an ending is correct using logic rules).
- **Minor Flaws**: The probability notation uses $\langle P(f(c, E) = e^+) \rangle$, but $\langle f \rangle$ is deterministic; it should clarify if $\langle f \rangle$ is probabilistic. $\langle N \rangle$ is assumed fixed (typically 4 in HellaSwag), but not specified.

Improved Symbolic Logic

Below is an improved version. I've expanded it to be more rigorous and comprehensive, merging it with elements from the mathematical framework (e.g., minimax, convergence) and generic framework (e.g., Goldilocks Zone). Key enhancements:

- **Formal Predicates for Commonsense**: Introduced symbolic predicates for reasoning types (e.g., physical plausibility, temporal consistency), enabling logical composition.
- **Logical Semantics**: Used modal logic (for "possible next events") and probabilistic elements to model plausibility.
- **Symmetric Handling of Endings**: Explicit rules for correct vs. adversarial negatives.
- **Complexity and Lattice**: Added theorem on complexity and a lattice for reasoning dimensions, with an algorithm for traversal.
- **Integration**: Tied to uncertainty (from Section 7), causality (from Section 12), and AF process.
- **Practical Extensions**: Included hybrid neuro-symbolic potential and examples for clarity.
- **Conciseness**: Removed redundancy, ensured mathematical consistency (e.g., probabilistic $\langle f \rangle$).

(3) Symbolic Logic - Improved

The HellaSwag task is formalized symbolically to capture commonsense reasoning in next-event prediction, bridging empirical NLP with formal logic. We model contexts as event sequences and endings as propositions, using modal logic for plausibility and predicates for reasoning types.

Definition 11 (Formal Logical Semantics - Improved):

Let $\langle \sigma(c) \rangle$ map context $\langle c \rangle$ (a sequence of events) to a set of logical formulas in modal logic, representing the situational state (e.g., using \Box for necessity, \Diamond for possibility). Let $\langle \pi(e) \rangle$ map ending $\langle e \rangle$ to a proposition.

The correct ending $\langle e^+ \rangle$ satisfies:

```
\langle \text{HellaSwag}(c, E) = e^+ \rangle \iff \exists \pi. \langle \pi \rangle = \langle \pi(e^+) \rangle \wedge \sigma(c) \models \Diamond \pi \wedge \bigwedge_{e^- \in E^-} \neg \Diamond \pi(e^-)
```

where $\langle \models \rangle$ is modal entailment, $\langle E = \{e^+\} \cup E^- \rangle$, and negatives fail plausibility (e.g., violate commonsense). For probabilistic models, extend to: $\langle P(\Diamond \pi \mid \sigma(c)) > \theta \rangle$ for $\langle e^+ \rangle$, $\langle P(\Diamond \pi \mid \sigma(c)) < \theta \rangle$ for negatives, tying to uncertainty in Definition 8.

Example: For context $\langle c = \rangle$ "A person is chopping wood.", $\langle \sigma(c) = \langle \text{Holding(Axe)}, \text{Near(Wood)} \rangle \rangle$. Correct $\langle e^+ = \rangle$ "The wood splits." satisfies $\langle \Diamond \text{Splits(Wood)} \rangle$; negative $\langle e^- = \rangle$ "The axe flies away." satisfies $\langle \neg \Diamond \text{Flies(Axe)} \rangle$ (physical violation).

Definition 12 (Symbolic Taxonomy of Commonsense Reasoning):

Reasoning types are predicates over $\langle (\sigma(c), \pi(e)) \rangle$, with external knowledge $\langle K \rangle$ (e.g., scripts, physics). The ending is correct if: $\langle \bigvee_r r(\sigma(c), \pi(e)) \rangle$, and adversarial if superficially true but $\langle \bigvee_r \neg r(\sigma(c), \pi(e)) \rangle$.

- **Physical Plausibility** ($\langle (\text{Phys}(\sigma, \pi)) \rangle$): $\langle \sigma(c) \cup K \models \Diamond \pi \rangle$ (consistent with physics, e.g., gravity).
- **Temporal Consistency** ($\langle (\text{Temp}(\sigma, \pi)) \rangle$): $\langle \sigma(c) \models \text{Next}(\pi) \rangle$ (logical sequence, using temporal operators).
- **Social Norms** ($\langle (\text{Soc}(\sigma, \pi)) \rangle$): $\langle \sigma(c) \cup K \not\models \text{bot} \text{ if } \pi \rangle$ (no social inconsistency).
- **Script Knowledge** ($\langle (\text{Script}(\sigma, \pi)) \rangle$): $\langle \exists s \in K. s(\sigma(c)) \implies \pi \rangle$ (matches event script, e.g., "cooking" implies "eating").
- **Causal Inference** ($\langle (\text{Caus}(\sigma, \pi)) \rangle$): $\langle \sigma(c) \models \pi \text{ causes } \text{Outcome} \rangle$ (from causal framework in Definition 13).

This integrates with Definition 4's reasoning-specific distributions: Each $\langle r \rangle$ corresponds to $\langle P(r \mid c, E) \rangle \approx 1/N$ for models but 1 for humans.

Theorem 6 (Complexity Classification - Improved):

Deciding if an ending is correct (is $\text{HellaSwag}(c, E) = e?$):

- PSPACE-complete for full modal logic (model checking with alternations).
- NP-complete for propositional fragments (satisfiability for plausibility).
- P for Horn-like scripts (forward chaining on rules).

Proof Sketch: Reduction from QBF for PSPACE; from 3-SAT for NP; efficient for definite clauses.

This aligns with Corollary 2's sample complexity: High VC-dimension requires large $\mathcal{O}(n)$ for generalization beyond biases.

Definition 13 (Reasoning Dimension Lattice - Improved):

Let $\mathcal{R} = \{\text{Phys, Temp, Soc, Script, Caus}\}$. Define lattice $\mathcal{L} = (2^{\mathcal{R}}, \subseteq, \sqcup, \sqcap)$:

- \sqcup (join): Union (e.g., Phys \sqcup Temp = physical-temporal reasoning).
- \sqcap (meet): Intersection (e.g., Script \sqcap Caus = script-based causality).

Orders by difficulty: Phys < Temp < Soc < Script < Caus (from basic to abstract). The Goldilocks Zone is the sublattice where model accuracy $\approx 1/|E|$ but human ≈ 1 .

Algorithm 4 (Lattice-Based Reasoning Traversal):

To symbolically compute $\text{HellaSwag}(c, E)$:

1. Compute $\sigma = \sigma(c), \{\pi(e) \text{ for } e \in E\}$.
2. Start at bottom \perp (no reasoning).
3. For each $r \in \mathcal{R}$ (ascending lattice order):
 - If $r(\sigma, \pi(e))$ for some e : Update current = current $\sqcup \{r\}$; candidate = e .
4. If $|\text{candidates}| = 1$: Return e (correct).
5. If multiple/none: Compute uncertainty $u(c, E)$ from Algorithm 2; if $u > \epsilon$, mark adversarial.

This is $\mathcal{O}(|\mathcal{R}| * |E| * |\sigma|)$, suitable for symbolic verifiers, and hybridizable (e.g., neural models for σ extraction, logic for validation).

Integration with Framework:

This symbolic logic enhances the Discriminator D in the generic framework: D can use modal entailment for filtering. For AF minimax (Definition 1), adversarial negatives maximize violations in the lattice. Causal aspects (Theorem 9) allow interventions like $\text{do}(\text{next event})$. Provides explainability by tracing failed predicates.

This improved formulation transforms HellaSwag's symbolic logic into a robust, verifiable structure, supporting co-evolution of benchmarks and models while highlighting genuine commonsense gaps.

—

(4) New Generic Document: A Framework for Adversarial Benchmark Creation

Title: A Generic Framework for Adversarial Benchmarking in AI Evaluation

1. Introduction

This document outlines a generic framework for creating benchmarks that robustly evaluate AI capabilities by minimizing the exploitation of superficial biases. The core insight is that static benchmarks are quickly "solved" not by mastering the underlying task, but by learning dataset-specific artifacts. This framework proposes a dynamic, adversarial approach to benchmark creation.

2. Problem Formulation

The goal is to create a benchmark dataset \mathcal{B} for a task \mathcal{T} that is:

- * **Trivial for Domain Experts:** Human performance P_H is near-perfect.
- * **Challenging for State-of-the-Art AI Models:** Model performance P_M is significantly lower than P_H , even after training on \mathcal{B} itself.
- * **Robust:** The difficulty stems from the core challenges of task \mathcal{T} , not from out-of-distribution trickery. The test distribution should be similar to the train distribution.

3. Core Framework: The Adversarial Co-Evolution Loop

The benchmark is built via an iterative loop involving three components: a **Generator**, a **Discriminator**, and a **Human Validator**.

Component A: The Generator (\mathcal{G})

- * **Role:** Produces challenging negative examples (incorrect answers/inputs) for the benchmark.
- * **Requirements:** Must be a state-of-the-art model for the domain (e.g., a large language model for text tasks). It should be capable of producing outputs that are superficially plausible but ultimately incorrect based on the rules of task \mathcal{T} .
- * **Process:** For a given seed (e.g., a context or a prompt), \mathcal{G} produces a large pool of candidate negative examples.

Component B: The Discriminator (\mathcal{D})

- * **Role:** Acts as a proxy for the current state-of-the-art. It is used to identify and filter out "easy" negative examples.
- * **Requirements:** Must be a powerful, representative model of the current best-performing architecture for the task.
- * **Process (Adversarial Filtering):**

1. A initial dataset is created with correct and generated negative examples.
2. **Iterate:** \mathcal{D} is trained on a subset of the data and used to evaluate a held-out set.
3. **Adversarial Replacement:** Any negative example in the held-out set that \mathcal{D} can easily identify as incorrect is replaced with a more challenging candidate from \mathcal{G} 's pool.
4. This loop continues until \mathcal{D} 's performance on the benchmark converges to a low level, indicating the negatives are highly adversarial.

Component C: The Human Validator (\mathcal{H})

- * **Role:** Ensures the benchmark's quality and validity.
- * **Process:** Humans review the final candidate examples, removing any where:
 - * The "correct" answer is ambiguous or wrong.
 - * A "negative" example is actually plausible or correct.

This ensures the benchmark's foundation is sound and that the difficulty for machines is not due to errors, but genuine task complexity.

4. The "Goldilocks Zone" of Complexity

A key finding is the existence of a "Goldilocks Zone" for the data. The examples must be:

- * **Complex Enough:** So that the generator \mathcal{G} makes believable mistakes that violate the underlying task logic (e.g., commonsense, physical laws).
- * **Simple Enough:** So that the discriminator \mathcal{D} struggles to detect these mistakes using surface-level patterns alone.

Finding this zone is empirical and crucial for the benchmark's success.

5. Evaluation and Iteration

- * **Benchmark Release:** The final dataset \mathcal{B} is released.
- * **Model Progress:** As new models (\mathcal{M}_{new}) are developed and achieve high performance on \mathcal{B} , the cycle restarts.
- * **Benchmark Evolution:** A new version of the benchmark (\mathcal{B}_{v2}) is created using \mathcal{M}_{new} as the new discriminator \mathcal{D} and the latest model as the generator \mathcal{G} . This ensures the benchmark continuously presents a frontier challenge.

6. Key Principles and Future Directions

- * **Principle of Co-Evolution:** Benchmarks should not be static; they must evolve with model capabilities.
- * **Principle of Adversarial Construction:** The best way to find a model's weaknesses is to have it fight against itself during dataset creation.
- * **Focus on Generalization:** Include "zero-shot" splits where the concepts being tested are unseen during training to measure true generalization versus memorization.
- * **Beyond Language:** This framework can be applied to other domains like vision (e.g., adversarial image generation) and robotics (e.g., challenging simulation scenarios).

This framework provides a blueprint for creating benchmarks that serve as true north stars for AI progress, ensuring that measured improvements reflect genuine advances in reasoning and understanding rather than shortcut learning.