# PoliScore:
# A Framework and Benchmark Suite for Policy Quality Engineering

Version 0.1 — December 1, 2025

*Disclaimer: The methodologies, frameworks, and evaluation systems described in this white paper represent an ongoing research and engineering effort. Certain components are partially implemented or in active development. Current production features of PoliScore may not yet reflect the full framework outlined herein.*

**Abstract**

This white paper introduces PoliScore, a first-principles framework for evaluating the structural quality of public policy, and PoliBench, a companion benchmark suite for assessing AI systems' competence in policy reasoning. Despite the central role of legislation in shaping societal outcomes, there is currently no unified, non-partisan standard for determining whether a proposed policy is well-designed, feasible, fiscally coherent, institutionally sound, and likely to achieve its stated goals. PoliScore addresses this gap by defining seven core dimensions of policy quality—problem clarity, evidence support, implementation feasibility, economic sustainability, distributional impact, governance integrity, and systemic risk—derived from foundational insights in political philosophy, welfare economics, institutional theory, and systems engineering.

PoliBench operationalizes this framework into reproducible evaluation tasks that test an AI system's ability to identify causal mechanisms, anticipate unintended consequences, detect governance vulnerabilities, analyze distributional effects, and evaluate feasibility constraints directly from legislative text. Together, PoliScore and PoliBench provide the foundation for a new discipline of policy quality engineering, offering structured, transparent, and interpretable methods for analyzing public policy prior to implementation. The paper concludes by outlining limitations, ethical considerations, the role of constrained web search in evidence retrieval, and opportunities for future research and institutional collaboration.

# Contents

# 1 Introduction

Public policy shapes the daily lives of individuals and the long-term trajectory of nations. Yet despite its immense influence, there is no widely accepted, non-partisan standard for evaluating the quality of legislation before it is enacted. Policymakers, analysts, and the public rely on fragmented tools such as budget scoring, partisan think-tank reports, post-hoc program evaluations, or ideological frameworks. These tools, while useful in isolation, do not provide a comprehensive or predictive understanding of whether a proposed policy is well-designed, feasible, equitable, or beneficial in practice.

Advances in artificial intelligence now make it possible to analyze complex policy texts at scale, but AI alone cannot fill this gap unless grounded in a coherent philosophical and methodological foundation. What is needed is not merely an automated analysis tool, but a rigorous theory of policy quality capable of guiding, constraining, and evaluating both human and machine reasoning about legislation.

PoliScore introduces such a framework, supported by PoliBench, a benchmark suite that operationalizes these ideas into concrete, reproducible evaluation tasks.

# 2 Foundational Principles of Policy Quality

## 2.1 Human Needs as the Basis of Societal Outcomes

At the core of any political system lies a simple and empirically grounded truth: human beings have universal, predictable needs. These include, at minimum, the requirements for survival (food, water, shelter, health, safety) and the conditions for flourishing (education, opportunity, autonomy, stability, and participation in society).

This principle is not ideological. It traces broadly through the history of philosophy and public policy, from classical notions of basic welfare to modern frameworks such as:

- Maslow's hierarchy of needs (Maslow, 1943)

- Sen's capabilities approach (Sen, 1999)

- Rawlsian justice as fairness (Rawls, 1971)

- Nussbaum's list of central human capabilities (Nussbaum, 2006)

- The UN Human Development Index (UNDP, 1990)

- OECD Better Life Index (OECD, 2011)

The implication is straightforward: policies ultimately exist to alter societal conditions in ways that affect these human needs. Therefore, any concept of policy quality must begin with the recognition of these universal determinants of human well-being.

## 2.2  The Purpose of Public Policy: Maximizing Societal Benefit

If human needs are universal, then public policy can be defined as:

> *A structured intervention intended to alter social, economic, or political conditions to maximize societal well-being, minimize harm, and ensure long-term stability.*

This definition integrates insights from multiple traditions:

- Utilitarianism — maximizing aggregate well-being.

- Rawlsian justice — ensuring fairness and protection for the least advantaged.

- Capabilities theory — expanding real freedoms and opportunities.

- Institutional economics — ensuring efficiency, stability, and low transaction costs (North, 1990).

- Governance theory — designing institutions that are transparent, accountable, and resilient (Ostrom, 1990).

From these traditions emerges a balanced, non-ideological goal:

> *Good policy is that which improves human outcomes without creating disproportionate harm, fragility, inequality, or institutional dysfunction.*

This is the foundational goal on which the PoliScore evaluative framework is built.

## 2.3  From First Principles to a Standard Rubric of Policy Quality

If

(a) humans have predictable needs, and

(b) the purpose of policy is to maximize societal benefit while minimizing harm,

then it follows that policy quality can be systematically evaluated according to how effectively a proposal moves society toward those outcomes.

PoliScore formalizes this into seven universal dimensions derived directly from these foundational principles:

1. **Problem Clarity & Causal Validity**
   Does the policy accurately diagnose the underlying issue and target the relevant causal mechanisms?
   *Grounding: Popper (1957); Sen (1987).*

2. **Evidence Base & Empirical Support**
   Is the proposed intervention supported by empirical research, historical precedent, or meaningful comparative data?
   *Grounding: Evidence-based policy literature.*

3. **Implementation Feasibility**
   Can existing institutions realistically execute the policy given resource, logistical, administrative, and temporal constraints?
   *Grounding: Herbert Simon (1947); implementation science.*

4. **Economic Efficiency & Fiscal Sustainability**
   Does the policy use resources responsibly, minimize waste, and avoid unsustainable long-term obligations?
   *Grounding: North (1990); welfare economics.*

5. **Distributional Impact & Fairness**
   How are benefits and burdens distributed across populations, and does the policy unjustifiably disadvantage certain groups?
   *Grounding: Rawls (1971); Sen (1999).*

6. **Governance Integrity & Institutional Risk**
   Does the policy maintain transparency, accountability, and resilience while minimizing opportunities for corruption or abuse?
   *Grounding: Ostrom (1990); public choice theory.*

7. **Unintended Consequences & Systemic Risk**
   Does the policy introduce fragility, perverse incentives, or cascading failures that undermine the intended outcomes?
   *Grounding: Popper (1945); Buchanan (1962).*

These seven pillars are not ideological criteria. They are derived from the intersection of:

- moral philosophy,

- economics,

- development theory,

- organizational behavior,

- risk analysis,

- governance studies,

- institutional design.

Taken together, they form a non-partisan, foundational theory of policy quality designed explicitly for computational evaluation.

## 2.4   Why This Framework Is Necessary

Existing institutions (such as the CBO, independent economic modeling groups, think tanks, and academic departments) evaluate policy through isolated lenses:

- cost impacts,

- economic forecasts,

- ideological alignment,

- advocacy positions,

- program evaluation after implementation.

None provide a unified, interdisciplinary, pre-implementation standard for determining whether legislation is well-designed, feasible, fair, evidence-based, and structurally beneficial.

PoliScore seeks to fill this gap by offering:

- a consistent methodology,

- grounded in decades of cross-disciplinary research,

- applicable directly to bill text,

- measurable, reproducible, and benchmarkable,

- independent of ideology.

This framework is the intellectual foundation for the PoliScore grading system and the PoliBench benchmark suite.

## 2.5   A New Field: Policy Quality Engineering

By synthesizing philosophical foundations, institutional economics, governance theory, and modern AI evaluation, PoliScore introduces what is effectively a new discipline:

> **Policy Quality Engineering:** the systematic, reproducible evaluation of public policy according to universal human needs, societal benefit, and institutional feasibility.

This white paper establishes the theoretical basis for that discipline.

# 3   The PoliScore Framework

The PoliScore framework provides a systematic, non-partisan method for evaluating the quality of public policy based on its expected real-world impact, feasibility, and alignment with universal human needs. Derived from the foundational principles articulated in Section 2, the framework operationalizes these concepts through seven core evaluation dimensions, each representing a necessary component of high-quality legislation.

Each dimension is designed to be:

- philosophically grounded,

- empirically motivated,

- institutionally relevant,

- computationally assessable, and

- applicable directly to legislative text.

Together, these dimensions form a comprehensive rubric for determining whether a policy is constructed in a way that maximizes societal benefit while minimizing harm, inefficiency, and unintended consequences.

## 3.1   Problem Clarity & Causal Validity

Effective policy begins with a clear, accurate understanding of the problem it seeks to address. Vague or misdiagnosed problems lead to interventions that fail to produce meaningful improvements or that target symptoms rather than causes.

A policy demonstrates clarity and causal validity when:

- the problem is explicitly defined and empirically measurable;

- the underlying causal mechanisms are identified;

- the proposed intervention plausibly affects those mechanisms;

- the theory of change is coherent and logically sound.

Policies that rest on untested assumptions, moral panic, or ideological narratives—rather than a valid causal model—score poorly on this dimension.

**Philosophical grounding:** Popper (1957), Sen (1987).
**Practical relevance:** Policy failures due to incorrect problem diagnosis.
**Computational challenge:** Extracting causal structures from bill text.

## 3.2   Evidence Base & Empirical Support

High-quality policy proposals demonstrate clear grounding in empirical research, comparative case studies, or validated theoretical frameworks. This dimension evaluates whether the intervention is supported by evidence that similar approaches have succeeded elsewhere, or whether its expected outcomes are consistent with existing knowledge.

Relevant criteria include:

- citation of empirical findings or documented precedents;

- alignment with established best practices in relevant fields;

- avoidance of claims contradicted by available data;

- transparency about uncertainty and knowledge gaps.

Policies lacking empirical grounding may rely on wishful thinking or unproven assumptions, increasing the risk of unintended harm.

**Philosophical grounding:** Evidence-based policy literature.
**Practical relevance:** Avoiding "policy theater" and ineffective reforms.
**Computational challenge:** Mapping bill provisions to known empirical findings.

## 3.3   Implementation Feasibility

Even a theoretically sound policy can fail if it is infeasible to implement. Feasibility depends on the capacity of institutions, agencies, and local systems to carry out the policy's mandates with available resources, logistics, workforce, technology, and time.

This dimension evaluates:

- administrative complexity and burden;

- clarity of agency responsibilities;

- resource requirements (financial, human, technical);

- timeline realism;

- reliance on unavailable or overextended infrastructure;

- potential bottlenecks or bureaucratic overload.

Policies that impose unrealistic workloads, require nonexistent infrastructure, or centralize responsibility in ways that exceed institutional capacity receive low scores.

**Philosophical grounding:** Herbert Simon (1947), implementation science.
**Practical relevance:** Real-world failure modes of legislation.
**Computational challenge:** Inferring administrative burden from text structure.

## 3.4  Economic Efficiency & Fiscal Sustainability

Public policy must allocate resources responsibly, avoid generating structural inefficiencies, and maintain long-term fiscal viability. This dimension evaluates whether the benefits of the policy justify its costs, whether incentives are aligned with real-world behaviors, and whether it avoids unnecessary waste or economic distortion.

Considerations include:

- long-term funding stability;

- cost-benefit alignment;

- administrative overhead;

- market distortions or inefficiencies;

- externalities (positive or negative);

- dynamic economic effects and sustainability over time.

Policies that rely on implausible revenue assumptions, produce excessive deadweight loss, or generate persistent deficits score poorly.

**Philosophical grounding:** North (1990); welfare economics.
**Practical relevance:** Long-term macroeconomic stability.
**Computational challenge:** Mapping provisions to economic impacts.

## 3.5 Distributional Impact & Fairness

Public policies distribute benefits and burdens across different groups within society. A high-quality policy should not impose disproportionate harm on specific populations or create unjustifiable imbalances in who gains and who loses. This dimension evaluates how a policy's effects are spread across income levels, regions, industries, and demographic groups, and whether these effects align with broadly accepted principles of fairness and responsible governance.

Key considerations include:

- which groups receive the primary benefits of the policy;

- which groups bear the costs or risks;

- whether the distribution of impacts is reasonable and transparent;

- whether the policy inadvertently worsens existing disadvantages;

- whether the policy shifts burdens onto populations with limited capacity to absorb them.

This pillar does not require or assume equal outcomes. Instead, it assesses whether the distribution of impacts is justified, defensible, and consistent with the stated goals of the policy, and whether any imbalances introduce meaningful risk or undue harm.

**Philosophical grounding:** Rawls (1971); Sen (1999).
**Practical relevance:** Avoiding regressive or disproportionate policy effects.
**Computational challenge:** Inferring distributional impacts from textual provisions.

## 3.6 Governance Integrity & Institutional Risk

Policies exist within complex governance structures. This dimension evaluates whether a proposal strengthens or undermines institutional integrity, transparency, accountability, and the rule of law.

Key criteria:

- clarity of authority and decision-making processes;

- adequacy of oversight and accountability mechanisms;

- risks of corruption, abuse of power, or regulatory capture;

- concentration of unregulated authority;

- resilience to political manipulation;

- clarity in compliance requirements.

Policies that concentrate discretionary power in unaccountable agencies, lack oversight, or create opportunities for corruption receive lower scores.

**Philosophical grounding:** Ostrom (1990), Buchanan (1962).
**Practical relevance:** Preventing governance failures and political fragility.
**Computational challenge:** Detecting governance structures from text.

## 3.7   Unintended Consequences & Systemic Risk

Complex systems often respond unpredictably to policy interventions. This dimension assesses the extent to which a policy may produce harmful unintended consequences, including perverse incentives, moral hazard, market failures, bureaucratic overload, or cascading systemic risks.

Evaluative criteria include:

- creation of fragile dependencies;

- incentive misalignment;

- spillovers into adjacent systems;

- risk of black markets or evasion;

- increased systemic fragility or bottlenecks;

- insufficient fail-safes or fallback mechanisms.

Policies that appear beneficial in theory but introduce hidden structural costs or vulnerabilities receive lower scores.

**Philosophical grounding:** Popper (1945); systems theory.
**Practical relevance:** Preventing policy backfires and crises.
**Computational challenge:** Predicting multi-system interactions.

## Summary of the Framework

Together, these seven dimensions form a holistic, first-principles model of policy quality. They address not only a policy's intent and potential benefits, but also its feasibility, fairness, evidence base, institutional risks, and long-term sustainability.

The PoliScore framework is designed to:

- provide structured, rational evaluation of legislation;

- enable reproducible scoring across policies and time;

- support AI-assisted legislative analysis;

- inform policymakers, researchers, and the public;

- reduce reliance on ideological or partisan heuristics.

This framework also provides the theoretical foundation for PoliBench, the benchmark suite introduced in Section 4, which tests whether AI systems can reliably interpret and evaluate policy quality along these seven dimensions.

# 4   The PoliBench Benchmark Suite

The PoliBench Benchmark Suite is a standardized set of tests designed to evaluate whether AI systems can accurately assess public policy along the seven dimensions of the PoliScore Framework. Whereas PoliScore provides the conceptual model for policy quality, PoliBench operationalizes that model into concrete, reproducible tasks that measure an AI system's ability to reason about legislative intent, feasibility, consequences, and institutional design.

PoliBench is not intended as a performance leaderboard for general AI capabilities. Rather, it is a domain-specific benchmark focused on policy reasoning, causal inference, and institutional awareness—areas where existing language models often demonstrate gaps despite strong natural language proficiency. The benchmark enables systematic comparison across AI systems and provides an empirical foundation for evaluating progress in computational policy analysis.

## 4.1   Motivation

Despite rapid advances in large language models (LLMs), there is currently no standardized method for testing their ability to interpret legislation or assess public policy quality. Existing AI benchmarks measure skills such as:

- question answering (SQuAD, Natural Questions),

- general knowledge (MMLU),

- reasoning (GSM8K, ARC),

- truthfulness (TruthfulQA),

- code generation (HumanEval).

None of these capture the skills required for policy evaluation, such as:

- recognizing ambiguous or misleading problem statements;

- detecting infeasible mandates;

- identifying governance risks;

- reasoning about distributional impacts;

- understanding institutional constraints;

- anticipating unintended consequences;

- evaluating evidence claims in context.

Public policy is a systems problem involving economics, governance, logistics, human behavior, and institutional dynamics. PoliBench fills a critical gap by testing whether AI systems can navigate these complexities in a disciplined and consistent way.

## 4.2 Objectives

PoliBench is designed to achieve five key objectives:

(1) **Evaluate policy-specific reasoning**
Measure whether an AI system can analyze legislative text in ways aligned with the PoliScore pillars.

(2) **Provide reproducible, standardized tests**
Ensure that models are evaluated under identical conditions, enabling meaningful comparisons.

(3) **Identify structural weaknesses in AI policy analysis**
Pinpoint which dimensions (e.g., feasibility, unintended consequences) pose the greatest difficulty for current models.

(4) **Support model improvement**
Provide researchers with targeted diagnostics for training and fine-tuning AI systems on legislative reasoning tasks.

(5) **Promote transparency and accountability**
Allow policymakers, academics, and the public to understand the strengths and limitations of AI in this domain.

## 4.3 Benchmark Structure

PoliBench is organized into seven test suites, one for each dimension of the PoliScore Framework. Each suite contains multiple task types, designed to assess both conceptual understanding and applied reasoning.

## Suite 1: Problem Clarity & Causal Validity

Tasks include:

- identifying unclear or incorrectly diagnosed problems;

- matching policy language to causal mechanisms;

- detecting irrelevant or logically inconsistent interventions.

## Suite 2: Evidence Base & Empirical Support

Tasks include:

- identifying claims not supported by evidence;

- evaluating whether a policy is consistent with known best practices;

- recognizing when a bill cites evidence that does not support its claims.

## Suite 3: Implementation Feasibility

Tasks include:

- detecting unrealistic administrative burdens;

- evaluating resource requirements;

- identifying conflicts with existing infrastructure or legal frameworks.

## Suite 4: Economic Efficiency & Fiscal Sustainability

Tasks include:

- predicting distortive market effects;

- identifying unsustainable funding structures;

- assessing whether proposed benefits justify projected costs.

**Suite 5: Distributional Impact & Fairness**

Tasks include:

- identifying groups disproportionately helped or harmed;

- evaluating burden-shifting dynamics;

- detecting regressive or unjustifiable distribution patterns.

**Suite 6: Governance Integrity & Institutional Risk**

Tasks include:

- identifying power concentrations without oversight;

- detecting corruption or capture risks;

- evaluating institutional clarity and accountability mechanisms.

**Suite 7: Unintended Consequences & Systemic Risk**

Tasks include:

- detecting perverse incentives;

- forecasting cross-system impacts;

- anticipating fragile dependencies or cascading failures.

Each suite includes constructed counterexamples, fictional bills, paired policy comparisons, and scenario-based prompts that allow objective scoring.

## 4.4   Benchmark Format

Each PoliBench task is designed to be:

- *deterministic* — clear pass/fail or graded criteria;

- *grounded* — tied to a specific quality dimension;

- *text-based* — directly applicable to legislative language;

- *model-agnostic* — usable by any AI system;

- *transparent* — accompanied by human-annotated rationale and expected answer patterns.

Tasks follow one of five formats:

- multiple-choice reasoning tests,

- short-form explanation tasks,

- bill snippet analysis,

- policy comparison tasks,

- error detection tasks (spot-the-flaw).

Scoring is done through a combination of rubric-based evaluation and automated correctness checks.

## 4.5 Example Test Types

Below are illustrative examples from several suites.

**Feasibility Example.**
**Prompt:** A bill requires every rural county (population $< 5,000$) to operate a full-service emergency hospital within one year.
**Question:** Identify the primary implementation challenge.
**Expected:** Workforce shortages; infrastructure constraints; unrealistic timeline.

**Governance Risk Example.**
**Prompt:** A bill grants an agency director unilateral authority to allocate funds "as they see fit," without reporting requirements.
**Expected:** Lack of oversight; corruption risk; unclear accountability.

**Distributional Impact Example.**
**Prompt:** A tax credit is available only to households with mortgage interest payments.
**Expected:** Benefits flow to homeowners; renters receive no support; regressive impact.

(These examples will appear in the full benchmark documentation.)

## 4.6 Scoring and Evaluation

Each model receives:

- a per-suite score (0–100),

- a composite PoliBench score,

- diagnostic reports identifying which dimensions require improvement.

Scores reflect:

- correctness,

- reasoning quality,

- internal consistency,

- robustness across variations,

- sensitivity to subtle policy flaws.

By aggregating results across suites, PoliBench reveals the policy reasoning profile of an AI model.

## 4.7 Availability and Reproducibility

PoliBench is published as:

- an open dataset,

- a reproducible evaluation script,

- documentation for scoring and reasoning analysis,

- example model outputs,

- version-controlled benchmark updates.

This ensures transparency and supports peer replication, academic review, and future model development.

**Summary**

PoliBench operationalizes the PoliScore Framework into a rigorous test suite for evaluating AI policy reasoning. It provides the technical foundation for comparing models, understanding their limitations, and improving automated legislative analysis. By bringing structure, transparency, and reproducibility to this domain, PoliBench aims to establish a new standard for computational policymaking research.

# 5 Methodology

The methodology underlying PoliScore is designed to translate legislative text into a structured evaluation across the seven dimensions of policy quality defined in Section 3. The process combines computational analysis, domain-specific prompt engineering, and rubric-based scoring to produce a transparent, interpretable assessment of a policy proposal's strengths and weaknesses.

PoliScore does not replace political judgment or override democratic deliberation. Rather, it provides a systematic framework for analyzing legislation in a manner that is consistent, reproducible, and grounded in well-established principles of public policy, institutional design, economics, and governance.

The methodological approach consists of four integrated components:

1. Text Preparation,

2. Dimension-Level Evaluation,

3. Scoring and Aggregation,

4. Interpretability and Justification.

Each is described below.

## 5.1 Text Preparation

Before analysis, the legislative proposal undergoes standardized preprocessing to ensure that the AI model receives input in a structured and interpretable form.

### 5.1.1 Document Segmentation

Legislation is segmented into logical units, such as:

- sections,

- subsections,

- enumerated provisions,

- definitions,

- mandates,

- authorizations,

- appropriations.

This segmentation allows the model to analyze each component independently and in context.

### 5.1.2   Contextual Metadata

Where available, contextual information is included:

- policy domain (healthcare, tax policy, infrastructure, etc.),

- jurisdiction,

- sponsoring entity,

- historical or comparative precedents,

- relevant statutory references.

Metadata helps the model identify potential feasibility constraints or institutional conflicts.

### 5.1.3   Legislative Normalization

Formatting inconsistencies, redundant language, or irrelevant boilerplate are removed to reduce noise and improve clarity.

## 5.2   Dimension-Level Evaluation

For each of the seven pillars of the PoliScore Framework, PoliScore applies a structured, domain-specific evaluation prompt. Each prompt is designed to assess a specific dimension using criteria derived directly from the theoretical foundations in Sections 2 and 3.

Each dimension evaluation includes three components.

### 5.2.1 Prompted Analysis

The model is asked a structured set of questions—tailored to that dimension—to elicit reasoning about:

- causal mechanisms,
- feasibility constraints,
- economic impacts,
- distributional patterns,
- governance structures,
- potential unintended consequences,
- evidence or lack thereof.

These questions are standardized to ensure consistency across evaluations.

### 5.2.2 Rubric-Guided Scoring

The model's response is evaluated according to:

- accuracy,
- recognition of key issues,
- reasoning quality,
- internal consistency,
- coverage of required criteria.

Each dimension is scored on a 0–100 scale.

### 5.2.3 Cross-Checks and Internal Consistency Tests

PoliScore performs additional diagnostics, such as:

- asking alternative formulations of the same question,
- checking for contradictions across dimension outputs,

- verifying that feasibility, economic impacts, and unintended consequences align with one another.

Internal contradictions reduce the dimension score.

## 5.3   Scoring and Aggregation

Once all seven dimensions have been evaluated and scored, PoliScore aggregates the results into:

- dimension-level scores (0–100),

- a composite policy quality score (also 0–100).

### 5.3.1   Weighting Scheme

By default, all dimensions are equally weighted, reflecting their equal importance in determining overall policy quality. However, the weighting system is configurable to support:

- domain-specific weighting (e.g., feasibility may be weighted more highly for emergency-response legislation),

- researcher-defined weighting,

- sensitivity analysis.

### 5.3.2   Composite Score Interpretation

The composite PoliScore is not a moral or ideological rating. It represents a structural assessment of:

- internal coherence,

- feasibility,

- fairness,

- evidence alignment,

- economic sustainability,

- governance soundness,

- systemic risk profile.

Scores may be interpreted as follows:

$$80\text{--}100 : \text{Strongly constructed policy with minor weaknesses,}$$
$$60\text{--}79 : \text{Moderately sound policy with identifiable risks or gaps,}$$
$$40\text{--}59 : \text{Weak policy likely to face implementation or outcome challenges,}$$
$$< 40 : \text{Structurally unsound policy with serious risks or flaws.}$$

These ranges are descriptive, not prescriptive.

## 5.4 Interpretability and Justification

PoliScore emphasizes transparency and interpretability to ensure users understand why a policy received its score.

### 5.4.1 Dimension Summaries

For each pillar, PoliScore generates:

- a concise explanation of strengths,

- a list of identified weaknesses,

- references to specific legislative sections,

- a justification for the assigned score.

### 5.4.2 Cross-Dimension Insights

PoliScore highlights interactions across dimensions, such as:

- economic impacts that undermine feasibility,

- governance risks that create unintended consequences,

- distributional effects that contradict the stated purpose of the bill.

This supports holistic understanding.

### 5.4.3 Explanation Consistency Checks

To improve trustworthiness, PoliScore includes:

- redundancy tests (asking the model to justify scores multiple ways),

- adversarial prompts,

- contradiction detection.

Inconsistent or unreliable explanations are flagged.

## 5.5 Alignment With PoliBench

PoliScore's scoring process is tightly coupled to the PoliBench benchmark suite. Models deployed within the PoliScore pipeline are required to demonstrate proficiency on PoliBench prior to being used for real-world legislative evaluation. This ensures:

- baseline competence,

- reproducible reasoning quality,

- calibrated performance across dimensions,

- transparency about limitations.

PoliBench acts as both a qualifying exam and a diagnostic tool for improving the underlying model.

## Summary

The PoliScore methodology provides a rigorous, structured, and interpretable approach to evaluating policy quality. By combining legislative segmentation, dimension-specific evaluation, rubric-based scoring, and cross-dimension analysis, PoliScore offers a transparent and reproducible system for assessing the strengths and weaknesses of proposed legislation.

This section outlines the formal scoring process; Section 6 examines how this process compares to the approaches used by existing institutions such as the CBO, think tanks, and international policy organizations.

# 6 Comparison to Existing Institutions

Public policy evaluation is not a new endeavor. Governments, academic centers, think tanks, and international organizations have long attempted to analyze the consequences of legislation. However, their approaches are fragmented, domain-specific, and often limited to economic forecasting, post-hoc program evaluation, or ideologically motivated analysis.

This section compares PoliScore to the institutions most commonly involved in policy assessment, and clarifies how the framework complements existing tools while filling an unmet need in pre-implementation policy quality evaluation.

## 6.1 Congressional Budget Office (CBO)

The Congressional Budget Office provides cost estimates and economic projections for federal legislation. Its analyses are widely respected and intentionally nonpartisan. However, by statutory mandate, the CBO:

- does not evaluate whether a policy is fair, feasible, or well-designed;

- does not assess governance risks or institutional fragility;

- does not analyze unintended consequences outside fiscal or macroeconomic domains;

- does not provide judgments about whether a policy is "good" or "poor";

- focuses almost exclusively on budgetary impacts, not policy quality.

**How PoliScore differs.**  PoliScore:

- evaluates seven dimensions of policy quality rather than a single economic dimension;

- focuses on pre-implementation design soundness;

- assesses governance structure, feasibility, fairness, and systemic risk, which CBO does not;

- does not produce fiscal projections, but instead evaluates whether funding mechanisms are structurally coherent.

CBO answers: *"What will this cost?"*
PoliScore answers: *"Is this policy structurally sound, feasible, and beneficial?"*

## 6.2 Think Tanks (Brookings, Heritage, AEI, CAP, Cato, etc.)

Think tanks play a major role in shaping public discourse about policy. They often produce:

- whitepapers,
- policy briefs,
- economic analyses,
- advocacy reports,
- commentary.

However, nearly all think tanks—no matter their orientation—produce work shaped by underlying ideological commitments or donor priorities. As a result:

- evaluations differ radically across institutions;
- frameworks for analysis vary widely;
- there is no unified or nonpartisan definition of policy quality;
- methodology is often opaque or narrative-driven;
- conclusions may be advocacy-oriented rather than diagnostic.

**How PoliScore differs.**    PoliScore:

- does not advocate for policy positions;
- uses a transparent, standardized framework;
- applies the same evaluative criteria to all legislation;
- grounds analysis in political philosophy, institutional economics, and governance theory rather than ideology;
- produces structured, replicable outputs, not narrative persuasion.

Think tanks answer: *"Is this policy aligned with our values or goals?"*
PoliScore answers: *"Is this policy well-designed according to universal structural criteria?"*

## 6.3 Academic Public Policy and Political Science Programs

Academic programs teach:

- cost-benefit analysis,

- program evaluation,

- ethics and justice,

- public administration,

- implementation theory.

However:

- methods vary by institution and professor;

- frameworks are conceptual, not standardized;

- academics rarely evaluate policy before implementation;

- analyses are usually qualitative, not rubric-based;

- no unifying cross-disciplinary "policy quality standard" exists.

**How PoliScore differs.**   PoliScore:

- operationalizes academic theory into a single coherent rubric;

- formalizes seven dimensions into measurable criteria;

- evaluates policy pre-implementation, not only after programs fail or succeed;

- integrates insights from economics, philosophy, governance, and systems engineering.

Academia answers: *"How should we think about policy?"*
PoliScore answers: *"How well-constructed is this specific policy?"*

## 6.4 International Organizations (OECD, UNDP, WHO, World Bank)

These institutions evaluate:

- development programs,

- governance indicators,

- effectiveness of existing policies,

- public health interventions,

- economic outcomes.

Their analyses are valuable but limited in scope. They mostly:

- evaluate implemented programs, not proposed legislation;

- rely on national data and long-term outcomes;

- focus on specific domains (health, governance, development);

- use retrospective evaluation rather than predictive structural analysis.

**How PoliScore differs.**   PoliScore:

- evaluates legislation prospectively, before implementation;

- operates at the level of bill text, not national indicators;

- applies the same framework across all policy domains;

- focuses on structural soundness, not outcome measurement.

International organizations answer: *"Did this policy improve development outcomes?"*
PoliScore answers: *"Is this policy likely to produce good outcomes?"*

## 6.5   Independent Economic Forecasting Models (Moody's, Penn Wharton Budget Model)

These models simulate macroeconomic consequences of policy proposals. They are sophisticated and data-intensive, but inherently narrow. They assess:

- GDP impact,

- employment effects,

- revenue outcomes,

- inflation or growth metrics.

They do not assess:

- governance risks,
- causal coherence,
- feasibility,
- distributional ethics,
- unintended consequences,
- structural soundness.

**How PoliScore differs.**   PoliScore:

- does not simulate macroeconomic variables;
- evaluates the design quality of the policy itself;
- includes governance, feasibility, and fairness—domains outside macro modeling.

Economic forecasters answer: *"What economic effects might this policy have?"*
PoliScore answers: *"Is this policy well-constructed across all relevant dimensions?"*

## 6.6   Summary of Differences

Table 1 summarizes the relationships between existing institutions and PoliScore.

## 6.7   Positioning

PoliScore does not replace existing institutions. It fills the missing analytical layer that sits between:

- academic theory,
- economic forecasting,
- practical governance,
- AI reasoning,
- public understanding.

By providing a unified, domain-agnostic framework for evaluating policy quality, PoliScore enables more structured public discourse, supports researchers, and creates new opportunities for interdisciplinary collaboration.

| Institution Type | What They Evaluate | What They Do Not Evaluate | PoliScore's Contribution |
|---|---|---|---|
| CBO | Budgetary/fiscal impacts | Governance, feasibility, fairness, systemic risk | Provides a holistic, pre-implementation structural evaluation. |
| Think Tanks | Ideology-driven arguments | Standardized, nonpartisan evaluation | Supplies neutral, structured scoring across seven dimensions. |
| Academia | Theory, ethics, evaluation methods | Unified applied rubric, operational scoring | Operationalizes academic theory into a consistent framework. |
| International Orgs | Retrospective outcomes | Pre-implementation analysis | Evaluates proposals before they are enacted. |
| Economic Models | Macro projections | Design quality, governance, distributional reasoning | Complements economic forecasts with structural analysis. |

Table 1: Comparison of existing institutions and PoliScore.

# 7 Limitations, Risks, and Ethical Considerations

While PoliScore and the PoliBench Benchmark Suite provide a structured and theoretically grounded framework for evaluating public policy, they also introduce methodological, computational, and ethical challenges. Recognizing these limitations is essential for responsible use and for ensuring that PoliScore complements, rather than replaces, democratic decision-making and expert judgment.

This section outlines the primary limitations of the framework, the risks associated with AI-assisted policy evaluation, and the ethical considerations necessary for its responsible deployment.

## 7.1 Limitations of the Framework

### 7.1.1 Incomplete Representations of Policy Context

Legislative text does not always contain:

- administrative history,

- political constraints,

- agency capabilities,

- cultural factors,

- stakeholder incentives,

- implementation environment.

PoliScore evaluates text as written, not the political or institutional context surrounding it. Real-world outcomes may differ from what the text alone suggests.

### 7.1.2 Dependence on Model Interpretation

PoliScore relies on the reasoning ability of large language models. While PoliBench validates baseline competence, no model is infallible. AI systems may:

- misinterpret ambiguous sections,

- overlook subtle governance issues,

- fail to identify complex incentive structures,

- inconsistently explain their reasoning,

- exhibit sensitivity to prompt phrasing.

Human oversight remains essential.

### 7.1.3 Normative Judgments Embedded in Pillar Definitions

Although the seven dimensions are derived from widely accepted principles, any evaluative framework contains implicit assumptions. For example:

- principles of "fairness" rely on philosophical traditions;

- feasibility depends on assumptions about institutional capacity;

- systemic risk depends on the evaluator's interpretation of fragility.

PoliScore mitigates this by grounding dimensions in well-established academic literature, but the framework is not value-free.

### 7.1.4 Challenges in Quantifying Qualitative Constructs

Some aspects of policy quality—such as institutional trust, political legitimacy, or cultural acceptance—are inherently difficult to quantify. PoliScore focuses on the structural quality of the policy text, but certain societal outcomes are fundamentally complex and unpredictable.

### 7.1.5 Early-Stage Field

Policy quality engineering is a new field. As such:

- the framework will evolve;

- dimensions may be refined;

- additional pillars may emerge;

- new benchmarks may be required as AI systems improve.

The methodology is intentionally designed to be iterative.

## 7.2 Risks Associated With AI-Assisted Policy Evaluation

### 7.2.1 Overreliance on AI Outputs

AI-generated evaluations, if misunderstood as authoritative, may:

- overshadow legitimate political debate;

- reduce the role of experts;

- disincentivize democratic deliberation;

- be mistaken for objective truth.

PoliScore is a tool for structured analysis, not a substitute for policymaking.

### 7.2.2 Risk of Misuse or Politicization

Any scoring system can be misused or selectively weaponized. Risks include:

- cherry-picking PoliScore results to support partisan narratives;

- misrepresenting composite scores without context;

- selectively highlighting favorable dimensions;

- using scores to target political opponents.

Mitigation requires transparency and publication of dimension-level details, not just composite numbers.

### 7.2.3   Model Bias and Blind Spots

Even when the framework is neutral, AI models may incorporate biases from:

- training data,

- institutional assumptions,

- cultural norms,

- market incentives,

- political rhetoric embedded in public discourse.

PoliBench mitigates this by testing for reasoning quality rather than surface-level pattern matching, but model bias cannot be fully eliminated.

### 7.2.4   Vulnerability to Adversarial Prompts

Sophisticated actors could attempt to influence model outputs by:

- strategically crafting bill text;

- embedding misleading language;

- exploiting LLM weaknesses;

- introducing ambiguous provisions.

Human review remains essential when evaluating high-stakes legislation.

## 7.3 Ethical Considerations

### 7.3.1 Transparency and Explainability

Users must understand:

- how PoliScore generates evaluations;
- what each dimension score means;
- how the model reached its conclusions.

PoliScore includes explanation requirements and cross-checks, but continued emphasis on transparency is critical.

### 7.3.2 Human Oversight and Democratic Authority

AI models and algorithmic scoring frameworks must not:

- replace elected representatives;
- undermine democratic decision-making;
- give undue authority to computational outputs.

PoliScore supports human deliberation; it does not supplant it.

### 7.3.3 Domain-Specific Sensitivity

Policies differ dramatically by sector (healthcare, taxation, environment, etc.). Ethical use of PoliScore requires:

- awareness of domain-specific complexities;
- consultation with subject matter experts;
- recognition that some impacts exceed textual analysis.

AI should augment domain experts, not replace them.

### 7.3.4 Inclusivity in Framework Development

As PoliScore evolves, its legitimacy depends on:

- peer review;

- collaboration with academics;

- multidisciplinary input;

- feedback from policymakers, economists, legal scholars, and civic groups.

An inclusive development process ensures that the framework does not reflect narrow assumptions or blind spots.

## 7.4   Mitigation Strategies

PoliScore incorporates several safeguards:

- PoliBench validation to ensure baseline reasoning quality;

- dimension-level transparency rather than opaque composite scores;

- cross-model comparison to detect inconsistencies;

- robust documentation for researchers;

- manual review options for analysts;

- versioning and changelogs to avoid unnoticed drift;

- public access to methodology to support scrutiny.

These mechanisms collectively reduce the risks associated with AI-assisted policy evaluation.

## Summary

PoliScore introduces a rigorous, first-principles approach to evaluating policy quality, but it is not a substitute for human judgment or democratic deliberation. Recognizing its limitations and potential risks is essential to its responsible use. By emphasizing transparency, interpretability, cross-checking, and academic grounding, PoliScore aims to serve as a constructive analytical tool rather than an authoritative arbiter of policy outcomes.

The next section, Section 8, outlines the integration of web search into PoliScore, the associated complexities, and the safeguards necessary for responsible use.

# 8 Integrating Web Search: Benefits, Complexities, and Concerns

As AI systems become increasingly capable of retrieving and synthesizing information from the web, responsible integration of external data into policy evaluation becomes both an opportunity and a challenge. Web search can significantly strengthen PoliScore's analytical depth, but it also introduces complexities related to reliability, bias, reproducibility, and governance. This section outlines the potential advantages of incorporating web search into the PoliScore methodology, along with the limitations and safeguards required to ensure responsible use.

## 8.1 Benefits of Web Search Integration

Web search provides access to real-world information that extends beyond the content of legislative text. When properly constrained, it can enhance accuracy and reduce reliance on incomplete model memory.

### 8.1.1 Access to Empirical Data

Many policy domains require up-to-date or domain-specific statistics, such as:

- workforce availability,

- infrastructure capacity,

- budget baselines,

- agency staffing levels,

- historical program enrollment,

- economic indicators.

These data points improve the quality of evaluations in:

- Pillar 2: Evidence Base & Empirical Support,

- Pillar 3: Implementation Feasibility,

- Pillar 4: Economic Sustainability.

### 8.1.2 Contextualizing Institutional Constraints

Web search can help identify:

- which agencies currently exist,

- what authority they have,

- notable oversight issues,

- past performance of similar programs,

- known bottlenecks in the administrative apparatus.

This supports:

- Pillar 6: Governance Integrity & Institutional Risk.

### 8.1.3 Reducing Hallucination

LLMs may hallucinate historical facts, legal structures, or administrative capacities when relying solely on internal training data. Web search:

- anchors evaluations in verifiable sources;

- reduces fabricated claims;

- improves factual grounding.

### 8.1.4 Enabling Policy Comparison

Web search allows the system to:

- retrieve prior legislation;

- compare cross-jurisdictional approaches;

- identify relevant case studies;

- verify whether similar policies have worked elsewhere.

This strengthens the evaluative process across multiple pillars.

## 8.2 Complexities and Technical Challenges

While web search offers clear benefits, its integration introduces several methodological and computational challenges.

### 8.2.1 Variability and Non-Reproducibility

The web changes continuously. Search results:

- vary by time;

- vary by location;

- may vary by query phrasing;

- may be influenced by search engine ranking algorithms.

This complicates reproducibility, a core requirement of scientific and academic evaluation frameworks.

### 8.2.2 Source Quality and Reliability

The open web contains:

- official data,

- academic research,

- advocacy reports,

- news articles,

- opinion pieces,

- misinformation,

- SEO-driven content.

Without strict filtering, AI may incorporate unreliable or biased information.

### 8.2.3 Fragmentation Across Policy Domains

Some domains (e.g., healthcare, taxation) have rich accessible data; others (e.g., cybersecurity, tribal governance, emerging technologies) may have sparse or inconsistent information. This asymmetry may lead to uneven evaluation quality across policy types.

### 8.2.4   Query Construction Ambiguity

Natural-language queries issued by an AI system may:

- be too broad or too narrow;

- return irrelevant content;

- introduce accidental bias;

- misinterpret data without proper grounding.

Careful prompt engineering and query constraints are required.

### 8.2.5   Latency and Performance

Real-time web search significantly increases:

- computational overhead,

- inference latency,

- system cost,

- user wait times.

This affects scalability for high-volume analysis.

## 8.3   Concerns and Risks of Web-Integrated Policy Evaluation

### 8.3.1   Bias Injection Through Search Engine Dynamics

Search engines prioritize:

- high-traffic sites,

- media outlets,

- pages optimized for engagement,

- politically charged topics.

This ranking dynamic may introduce hidden biases into the evaluation process, even if no partisan intent exists.

### 8.3.2 Political and Ideological Contamination

If search results include:

- think tank content,

- political commentary,

- advocacy messaging,

- editorial interpretations,

then model outputs may inherit the ideological predispositions of the sources.

### 8.3.3 Stability Over Time

As search engine algorithms evolve:

- rankings shift;

- domain authority changes;

- new narratives become dominant;

- old content is deprecated.

Evaluations may drift over time in ways unrelated to policy quality.

### 8.3.4 Difficulty Maintaining Neutrality

Without strict guardrails, web search could erode the non-partisan foundation of the PoliScore framework and undermine trust.

## 8.4 Safeguards and Responsible Integration Strategies

To mitigate risks, any web search integration should follow structured constraints.

### 8.4.1 Limited Source Categories

Restrict searches to:

- official government websites;
- academic institutions;
- nonpartisan statistical agencies;
- recognized international organizations;
- peer-reviewed research.

Exclude:

- opinion pieces,
- advocacy groups,
- partisan outlets,
- editorial content,
- think tank position papers.

### 8.4.2 Evidence-Type Filtering

Search should retrieve facts, not interpretations:

- numeric data,
- historical precedents,
- legal definitions,
- administrative roles,
- technical specifications.

### 8.4.3  Cached Snapshots for Reproducibility

To reduce time-based variability:

- store retrieved documents;

- version control search results;

- allow external auditing;

- include citations and timestamps.

### 8.4.4  Transparency in Usage

Reports should explicitly state:

- whether web search was used;

- which sources were cited;

- what data influenced the score;

- confidence levels based on evidence availability.

### 8.4.5  Human Oversight

Analysts and researchers should review flagged sections or cases where:

- search results contradict model expectations;

- data is sparse or ambiguous;

- governance risks are inferred from historical controversies.

AI augments human insight; it does not replace it.

## 8.5  Role of Web Search in the Future of Policy Quality Engineering

Web search can play a powerful but constrained role in strengthening policy evaluation frameworks. As AI models improve in retrieval accuracy and source discrimination, PoliScore may move toward:

- more robust evidence gathering,

- dynamic institutional context modeling,

- sophisticated counterfactual analysis,

- cross-jurisdictional comparisons.

However, maintaining neutrality and reproducibility will require:

- rigorous source control,

- clear methodological boundaries,

- transparent documentation,

- continuous monitoring for bias,

- strong collaboration with academic experts.

Web search should enhance the quality of evaluation—not overshadow or distort the structural analysis at the core of the PoliScore framework.

## Summary

Integrating web search into policy evaluation offers substantial advantages, including improved factual grounding, stronger feasibility analysis, and reduced hallucination. However, it introduces complexities related to source quality, bias, reproducibility, and governance. A careful, constrained, and transparent integration strategy allows PoliScore to benefit from real-world evidence while preserving its non-partisan, first-principles foundation.

# 9 Conclusion and Future Directions

Public policy shapes the foundational conditions under which individuals and societies live, work, and thrive. Yet until now, there has been no unified, non-partisan, and methodologically rigorous framework for evaluating the structural quality of legislation before it is enacted. PoliScore addresses this gap by offering a principled, interdisciplinary model grounded in human needs, political philosophy, institutional economics, governance theory, and systems thinking.

By formalizing policy quality into seven core dimensions—problem clarity, evidence support, feasibility, economic sustainability, distributional impact, governance integrity, and systemic risk—PoliScore introduces a reproducible standard for assessing the internal soundness of proposed legislation. The accompanying PoliBench Benchmark Suite operationalizes these principles into concrete tasks that evaluate whether AI systems possess the reasoning skills necessary to perform this analysis reliably.

Together, PoliScore and PoliBench represent early steps in what may become a broader field: policy quality engineering—a discipline focused on the structural, empirical, and institutional soundness of public policy design.

## 9.1 Contributions of This Work

This whitepaper establishes:

- a first-principles theory of policy quality grounded in universal human needs;

- a formal, seven-pillar framework for evaluating legislation;

- a benchmark suite (PoliBench) for measuring AI competence in policy reasoning;

- a methodology for generating transparent, reproducible, interpretable policy evaluations;

- a structured comparison to existing institutions and their limitations;

- an analysis of risks, ethical considerations, and safeguards involved in AI-assisted evaluation;

- a strategy for responsible integration of web search into the evaluation pipeline.

These contributions collectively create the foundation for a standardized, academically defensible approach to pre-implementation policy evaluation.

## 9.2 Opportunities for Future Research and Development

PoliScore and PoliBench are early prototypes of a much larger ecosystem that could emerge around computational policy evaluation. Several areas offer immediate opportunities for expansion.

### 9.2.1 Empirical Validation and Case Studies

Applying PoliScore to:

- historical legislation,

- failed policies,

- successful reforms,

- cross-national comparisons

will help validate the framework's predictive accuracy and refine its criteria.

### 9.2.2 Collaboration with Academic Institutions

Partnerships with universities, policy schools, and research centers can:

- stress-test the methodology;
- provide access to subject matter experts;
- ensure theoretical robustness;
- support peer review;
- enhance academic legitimacy.

### 9.2.3 Expansion of PoliBench

Future versions of the benchmark may include:

- domain-specific test suites (healthcare, taxation, climate policy, infrastructure);
- multi-step reasoning tasks;
- adversarial tests to detect bias or fragility;
- multilingual evaluations for international policy contexts.

### 9.2.4 Integration Into Civic and Governmental Processes

PoliScore could support:

- legislative drafting workflows;
- agency impact assessments;
- think tank analysis standards;
- civic education platforms;
- media fact-checking of policy claims;
- tools for voters seeking clear and structured information.

### 9.2.5 Enhanced Retrieval and Evidence Systems

Carefully constrained web search, specialized retrieval pipelines, and curated knowledge bases could strengthen:

- evidence evaluation;
- feasibility assessments;
- institutional context modeling;
- cross-jurisdictional comparisons.

### 9.2.6 Multimodal Policy Understanding

As models evolve, future work could incorporate:

- charts, budgets, and fiscal tables;
- legal diagrams and regulatory schemas;
- geospatial data;
- workflow diagrams for administrative processes.

## 9.3 A Path Toward a New Policy Evaluation Standard

PoliScore is not intended to replace political judgment, democratic deliberation, or human expertise. Rather, it aims to provide a clear, structured, transparent foundation for reasoning about legislation—one grounded in universal principles rather than ideology.

By elevating the discourse surrounding public policy from partisan narratives to structural analysis, PoliScore opens the door to:

- clearer public understanding,
- more responsible policymaking,
- better institutional design,
- more effective AI systems,
- a healthier democratic process.

The long-term vision is a world in which policymakers, researchers, journalists, and citizens alike have access to transparent, neutral, and rigorous assessments of policy quality—tools that help anchor public debate in reasoned analysis rather than rhetoric.

## 9.4 Closing Remarks

The challenges facing modern societies—economic transformation, climate change, technological disruption, public health, geopolitical instability—demand a new level of clarity and rigor in how legislation is evaluated. PoliScore represents a foundational step toward that future. By combining philosophical grounding, interdisciplinary rigor, and computational capability, it offers a framework for understanding not just what policy says, but how well it is constructed.

This whitepaper invites collaboration from academics, policymakers, technologists, and civic institutions to refine, validate, and expand this work. The development of policy quality engineering will require broad participation, diverse expertise, and ongoing scrutiny—but the potential benefits for democratic governance and societal well-being are substantial.

# References

[1]  John Rawls, *A Theory of Justice*. Harvard University Press, 1971.

[2]  Amartya Sen, *Development as Freedom*. Oxford University Press, 1999.

[3]  Elinor Ostrom, *Governing the Commons*. Cambridge University Press, 1990.