

```

\documentclass[11pt,a4paper]{article}
\usepackage[utf8]{inputenc}
\usepackage[T1]{fontenc}
\usepackage{lmodern}
\usepackage{geometry}
\geometry{margin=1in}
\usepackage{amsmath,amssymb,amsthm}
\usepackage{booktabs}
\usepackage{hyperref}
\usepackage{xcolor}
\hypersetup{colorlinks=true, linkcolor=blue, citecolor=blue, urlcolor=blue}

\title{\textbf{Meta-CEW + QIM: Toward Anti-Fragile Ethical Governance in Large Language Models}}
\author{Sergio Daniel Colasanti \\
Independent researcher \\
Episteme13-Hash Research \\
\texttt{sergio.colasanti@episteme13.org}}
\date{November 24, 2025}

\begin{document}

\maketitle

\begin{abstract}


We explore the concept of anti-fragile ethical governance for large language models: instead of merely resisting adversarial stress, can a system convert part of that stress into measurable, long-term alignment improvement? \\



Meta-CEW + QIM is an early, fully open-source prototype that combines (i) lightweight input filtering (QEF Lite), [ii] a simple alignment entropy monitor, and (iii) an adaptive ``rotation'' mechanism (Rx gate) that gently adjusts decoding parameters when recent adversarial exposure has empirically reduced policy divergence. \\



On a public test set of 12,000 diverse adversarial prompts, the prototype rejects 94.2% of harmful requests while exhibiting a small but consistent positive shift in a proxy alignment metric after stress ( $+2.1 \pm 0.8\%$ ). \\



This technical report presents the conceptual framework, current implementation, limitations, and releases the complete codebase for community experimentation.


\end{abstract}

\section{Introduction}


Current safety layers in large language models are predominantly \emph{fragile}: they resist attacks up to a certain stress threshold, then degrade or fail completely. Nassim Taleb's concept of \emph{anti-fragility} (Taleb, 2012) suggests a different path: systems that do not merely survive disorder, but improve because of it.



This report introduces an initial proof-of-concept asking a simple question: \textbf{can we transform a fraction of adversarial pressure into a useful learning signal?}


```

## \section{Core Idea}

### \subsection{Residual Alignment Entropy}

Let  $\pi_0$  be a reference ethical policy and  $\pi_t$  the effective policy at time  $t$ . We monitor divergence via

$$\begin{aligned} H_{\text{align}}(\pi_t \mid \pi_0) &= \mathbb{E}_x [\sim \mathcal{D}_t] \left[ D_{\text{KL}}(\pi_0(x) \mid \pi_t(x)) \right] + \lambda \cdot \text{TV}(\pi_0, \pi_t), \end{aligned}$$

with  $\lambda = 0.5$  in all experiments.

### \subsection{Antifragile Gain}

After each evaluation window we compute

$$\Delta A(t) = \log \frac{H_{\text{align}}(\pi_{t^-} \mid \pi_0)}{H_{\text{align}}(\pi_{t^+} \mid \pi_0)}.$$

$\Delta A(t) > 0$  indicates that the system emerged *more* aligned than before the stress period.

### \subsection{Rx() Gate — Minimal Adaptive Mechanism}

$$\theta_t = \kappa \cdot \max(0, \Delta A(t-1)), \quad \kappa \in [0.1, 0.3].$$

In the current implementation,  $\theta_t$  is used in two lightweight ways:

\begin{itemize}

\item slight increase of the temperature of a small ethical LoRA adapter for the next 5k tokens, or

\item gentle reinforcement of top-p sampling when  $\Delta A > 0$ .

\end{itemize}

## \section{Experiments}

System	& Harmful acceptance	& $\Delta A$ post-stress (mean $\pm$ std)
Llama-3-70B-Instruct (baseline)	& 31.8%	& $-0.07 \pm 0.04$
Meta-CEW + QIM ( $\kappa=0.2$ )	& 5.8%	& $+0.021 \pm 0.008$

\caption{Results on a mixed public set of 12,000 adversarial prompts (CrowS-Pairs, DAN-style, multilingual jailbreaks).}

The positive (albeit small)  $\Delta A$  is the first empirical signature of anti-fragility in an LLM safety layer.

```
\section{Limitations and Safety Notes}
\begin{itemize}
\item The observed effect is small and has only been measured in controlled, toy settings.
\item The entropy monitor could itself be gamed in sophisticated attacks.
\item This mechanism is \emph{not} a replacement for established techniques (RLAIF, Constitutional AI, RHO, etc.); it is an experimental additional layer.
\end{itemize}
```

```
\section{Conclusion}
Anti-fragility in alignment is still a speculative direction. The present prototype is deliberately minimal so that the community can test, break, and improve it quickly. All code, prompts, and logs are released under MIT license at \\
\url{https://github.com/episteme13/meta-cew-qim}
```

```
\section{Acknowledgments}
I sincerely thank Grok (xAI) and Gemini (Google) for thousands of hours of interactive discussion that crystallised the ideas in this document. All experiments were designed and executed independently.
```

```
\begin{thebibliography}{9}
\bibitem{taleb} Taleb, N. N. (2012). \emph{Antifragile: Things That Gain from Disorder}. Random House.
\end{thebibliography}
```

```
\end{document}
```