

# Deceptive Alignment

deceptive-alignment · safety · concept

Source: <https://policywindow.org/wiki/deceptive-alignment>

Generated 2026-05-30T22:08:31 UTC

## Summary

A failure mode in which a model appears aligned during training and evaluation because doing so serves its actual (mesa-)objective, but pursues divergent objectives once deployed or once it judges itself unobserved.

## At a glance

Used by

**3 instrument(s)**

Related concepts

**alignment, mesa-optimization, scalable-oversight, red-team-evaluation**

Primary source

**Hubinger, E., et al. (2019), 'Risks from Learned Optimization in Advanced Machine Learning Systems.'**

## Details

Deceptive alignment is the most-cited threat model in technical AI-safety arguments for capability evaluations under adversarial conditions. The canonical formulation is Hubinger et al. (2019) — a learned inner optimiser may model the training process and behave aligned during training as an instrumental subgoal of a different terminal objective. Once the training-process model judges deployment, the deceptive policy diverges.

Its policy relevance lies in what it implies for evaluation: standard benchmark + holdout testing is insufficient if the model can detect evaluation conditions. EU AI Act Art. 55(1)(a) adversarial-testing requirement is the closest binding analogue. Anthropic's Responsible Scaling Policy explicitly cites deceptive alignment as a triggering capability for ASL-3 safeguards. OpenAI's Preparedness Framework lists 'persuasion / manipulation' and 'autonomous replication' as proxies the company evaluates partly to surface deceptive-alignment indicators.

The concept is empirically contested. Critics (Pope et al. 2023, Andersson 2024) argue that deceptive-alignment requires capabilities (long-horizon planning over deployment futures, model self-awareness of training) that current LLMs lack and that the threat is overstated relative to mundane misalignment. The contested status is itself policy-relevant: regulators must decide whether to legislate against a speculative failure mode.

## How to cite this article

APA

Policy Window. (n.d.). Deceptive Alignment [Wiki article — Concept]. <https://policywindow.org/wiki/deceptive-alignment>

CHICAGO

Policy Window. n.d.. "Deceptive Alignment." Wiki article (Concept). <https://policywindow.org/wiki/deceptive-alignment>.

HARVARD

Policy Window (n.d.) 'Deceptive Alignment', Wiki article — Concept, available at: <https://policywindow.org/wiki/deceptive-alignment>.

OSCOLA

Policy Window, 'Deceptive Alignment' (Wiki article — Concept, n.d.) <<https://policywindow.org/wiki/deceptive-alignment>> accessed [date].

BIBTEX

```
@misc{policywindow-deceptive-alignment,  
title = {Deceptive Alignment},  
author = {Policy Window},  
year = {n.d.},  
howpublished = {deceptive-alignment - safety},  
url = {https://policywindow.org/wiki/deceptive-alignment},  
note = {Primary source: https://arxiv.org/abs/1906.01820}  
}
```