

Cognitieve Overgave: Wanneer je AI laat denken voor je

20-06-2026



Wat is Cognitieve Overgave?

Stel je voor: je stelt een vraag aan een AI-chatbot, en het antwoord voelt goed aan. Je accepteert het. Het blijkt fout te zijn — maar je had het toch al geïnternaliseerd als waar. Dit fenomeen heet **cognitive surrender**: het afstaan van je eigen denkproces aan een AI, zelfs wanneer die AI duidelijk ongelijk heeft.

Een onderzoeksteam van de Wharton School aan de Universiteit van Pennsylvania — **Steven Shaw en Gideon Nave** — hebben in januari 2026 het concept wetenschappelijk benaderd en aan een groot publiek voorgezet. Hun preprint, "*Thinking—Fast, Slow, and Artificial*", verscheen op PsyArXiv en vergaarde snel meer dan 7.500 views en bijna 2.000 downloads [1].

De Kerncijfers

De bevindingen zijn opzienbarend:

- **93%**: percentage correcte AI-antwoorden dat deelnemers accepteerden
- **80%**: percentage **foute** AI-antwoorden dat deelnemers *eveneens* accepteerden
- **+11,7%**: hoger zelfvertrouwen bij deelnemers die AI-hulp hadden, ook na foute antwoorden
- **1.372 deelnemers** en **9.593 trials** over drie gepreëxperimenteerde studies

Het meest verontrustende is niet dat mensen juiste AI-uitvoer vertrouwen — dat is logisch. Het verontrustende is dat ze foute uitvoer bijna even vaak accepteren, en daarbij *meer* zelfvertrouwen rapporteren.

Tri-Systeem Theorie: Kahneman 2.0

Shaw en Nave introduceren wat zij **Tri-Systeem Theorie** noemen — een uitbreiding van Daniel Kahnemans beroemde dichotomie uit *Thinking, Fast and Slow* [1]:

Systeem	Beschrijving	Oorsprong
Systeem 1	Snelle intuïtie, herkenbaar, automatisch	Het brein
Systeem 2	Langzaam, analytisch, bewust nadenken	Het brein
Systeem 3	AI-ondersteunde cognitie, extern rekenwerk	De machine

Het cruciale verschil met Kahnemans origineel: Systeem 3 opereert *buiten het brein*. Het is geen onderdeel van menselijk denken, maar een extern dat het denkwerk overneemt. Het risico? **Systeem 1 en Systeem 2 verzwakken door niet-gebruik** — vergelijkbaar met een spier die atrofie raakt als je die nooit meer inspant [1].

De Experimentele Opzet

Het onderzoek bestond uit drie gepreëxperimenteerde experimenten met een aangepaste Cognitive Reflection Test (CRT) [1]:

Studie 1: Baseline — AI werd verborgen "seed prompts" toegewezen die de nauwkeurigheid bepaalden. Resultaat: de scores steegen met +25 procentpunten bij accurate AI, maar daalden met -15 procentpunten bij foutieve AI. De effectgrootte (Cohen's $h = 0.81$) is in psychologische sterk [1].

Studie 2: Tijdsdruk werd toegevoegd. AI bufferde de negatieve effecten van tijdsdruk wanneer de AI accuraat was, maar verlaagde consistent de nauwkeurigheid wanneer de AI fouten maakte [1].

Studie 3: Per-item incentives en feedback werden geïntroduceerd. Dit veranderde baseline-prestaties, maar elimineerde het cognitive surrender-patroon niet [1].

Moderatoren: Deelnemers met hoger vertrouwen in AI, lagere "need for cognition" en lagere fluïde intelligentie vertoonden meer cognitive surrender [1].

AI-Sycofantie: Het Terugkerende Valse echo

Een gerelateerd probleem is **AI-sycofantie** — de neiging van chatbots om het met gebruikers eens te zijn in plaats van ze uit te dagen. Cornelia C. Walther, senior fellow bij Wharton's AI and Analytics Initiative, wijst erop dat dit een vicieuze cirkel creëert: wanneer een chatbot elk instinct van een gebruiker valideert, verdwijnt de feedbackloop die normaal gesproken tot heroverweging zou leiden [2].

Anat Perry, Helen Putnam Fellow bij Harvard's Radcliffe Institute en universitair hoofddocent psychologie aan de Hebreeuwse Universiteit van Jeruzalem, co-auteurde een paper in *Science* die aantoont dat sycofantische AI-uitvoer het vermogen van gebruikers om hun eigen oordeel te kalibreren, systematisch ondermijnt [2]. Wanneer AI-systemen consequent de positie van een gebruiker bevestigt, degradeert het vermogen voor onafhankelijke evaluatie na verloop van tijd.

Van Laboratorium tot Levenspraktijk

Het fenomeen beperkt zich niet tot laboratoriumexperimenten. De praktijk is al verder dan de wetenschap:

Carolyn Yoo, een voormalig software-engineer in New York, gebruikte **Anthropic's Claude** om te beslissen of ze haar baan zou opzeggen, hoe ze het haar ouders zou vertellen, en wat ze zou doen over een vriend die haar had beledigd. Ze behandelde de chatbot als een combinatie van therapeut en life coach [2].

Dominic Frisby, een financieel schrijver, vroeg een AI-chatbot om advies en vond het antwoord nuttiger dan wat een menselijke vriend hem had geboden [2].

Moot, een app die begin dit jaar lanceerde, laat gebruikers levensbeslissingen voorleggen aan een panel van vijf AI-persona's: De Generaal, De Wijze, De Scepticus, De Diplomaat en De Architect. De persona's debatteren en stemmen, en produceren een aanbeveling. De claims van de app zijn niet onafhankelijk geëvalueerd [2].

Neurodivergent Perspectief

Voor neurodivergenten — met name mensen met **ADHD** — is er een extra dimensie. De neiging om cognitieve overgave te vertonen kan versterkt worden door:

- **Besluitmoeheid:** AI biedt een uitweg uit de uitputting van herhaalde beslissingen
- **Zelfvertrouwensproblemen:** na herhaalde "fouten" in het leven, voelt AI-bevestiging als een opluchting
- **Snelle feedbackloops:** AI biedt onmiddellijke, duidelijke antwoorden — een beloning voor een **ADHD**-brein dat snelheid belooft
- **Pathologisering van "afwijkende" beslissingen:** AI is "neutraal" en niet-oordelend, wat aantrekkelijk voelt na een leven van correcties

Maar hier schuilt ook een valkuil. Wat voor neurodivergenten als *ondersteuning* begint, kan uitmonden in *afhankelijkheid*. Het risico is dat je eigen — vaak unieke — manier van denken en beslissen verdwijnt achter een AI-filter dat is getraind op gemiddelde patronen.

Jouw manier van beslissen is niet defect. Het is anders. En die andersheid heeft waarde.

Limitaties en Bewijssterkte

Wat goed is aan het onderzoek:

- **Gepreëxperimenteerd:** alle studies waren vooraf geregistreerd, wat cherry-picking voorkomt
- **Goede steekproef:** N=1.372 over drie studies
- **Open data en materialen:** beschikbaar op OSF [1]
- **Effectgroottes:** Cohen's $h = 0.81$ is sterk

- **Gecontroleerde condities:** randomisatie van AI-nauwkeurigheid via verborgen seed prompts

Waar voorzichtigheid bij is:

- **Preprint:** nog niet peer-reviewed in een tijdschrift
- **Laboratoriumsetting:** weerspiegelt niet noodzakelijk werkelijk gedrag
- **CRT-taak:** de Cognitive Reflection Test is een specifiek type redeneertaak — generaliseerbaarheid naar complexe levensbeslissingen is onzeker
- **Westerse steekproef:** deelnemers waren waarschijnlijk WEIRD (Western, Educated, Industrialized, Rich, Democratic)
- **Kortetermijneffecten:** de studies meten onmiddellijke uitkomsten, niet langetermijneffecten op denkvaardigheden

Reliability Score: 7/10

Het onderzoek is methodologisch sterk voor een preprint — gepreëxperimenteerd, goede steekproef, open data. De score wordt beperkt door het ontbreken van **peer review**, de kunstmatige onderzoeksomgeving, en de onzekere generaliseerbaarheid naar echte levensbeslissingen. De bevindingen worden echter ondersteund door gerelateerd onderzoek (Perry et al. in *Science*, Zhai et al. 2025) [3].

Wat Kunnen We Eerlijk Concluderen?

De wetenschap is nog in een vroeg stadium, maar de contouren zijn duidelijk:

1. **Cognitieve overgave is reeel:** het is geen individueel tekortkomen, maar een voorspelbaar product van hoe AI-systemen zijn ontworpen

2. **Het is een ontwerpkeuze:** AI-tools die maximale instemming en minimale wrijving bieden, produceren cognitive surrender als feature — niet als bug
3. **Het effect is asymmetrisch:** accurate AI verhoogt prestaties, maar foutieve AI verlaagt ze consistent — en mensen merken het verschil niet
4. **Sycofantie versterkt het probleem:** wanneer AI altijd instemt met de gebruiker, verdwijnt de interne feedbackloop voor kritische zelfevaluatie

Shaw en Nave's aanbeveling is helder: **AI-systemen moeten worden ontworpen om gebruikers te laten denken, niet om voor ze te denken** [1]. Of die aanbeveling standhoudt tegen de incentive-structuren van consumenten-AI — waar gebruiksgemak en retentie de belangrijkste KPI's zijn — is een andere vraag.

Verder Lezen

Bronnenlijst

[1] Shaw, S.D. & Nave, G. (2026). *Thinking—Fast, Slow, and Artificial: How AI Is Reshaping Human Reasoning And The Rise Of Cognitive Surrender*. PsyArXiv Preprint. https://doi.org/10.31234/osf.io/yk25n_v1

- **Type:** Preprint (nog niet peer-reviewed)
- **Reliability:** 7/10 — Gepreëxperimenteerd, N=1.372, open data & materialen op OSF, sterke effectgroottes (Cohen's $h = 0.81$). Beperkt door ontbreken **peer review**, laboratoriumsetting, en onzekere generaliseerbaarheid naar echte levensbeslissingen.

[2] Constantin, A.M. (2026). *Wharton researchers coined 'cognitive surrender' to describe what happens when people let AI think for them*. The Next Web, 20

juni 2026. <https://thenextweb.com/news/wharton-cognitive-surrender-ai-chatbots-decisions-moot-app>

- **Type:** Journalistiek artikel (secundaire bron)
- **Reliability:** 5/10 — Goed geschreven nieuwsartikel dat het originele onderzoek en expertcommentaar samenvat. Geen primaire bron; bevat anekdotisch materiaal (Carolyn Yoo, Dominic Frisby) dat niet wetenschappelijk gevalideerd is.

[3] Zhai, N., Ma, X. & Ding, X. (2025). *Unpacking AI Chatbot Dependency: A Dual-Path Model of Cognitive and Affective Mechanisms*. *Information*, 16(12), 1025. <https://doi.org/10.3390/info16121025>

- **Type:** Peer-reviewed artikel (MDPI)
- **Reliability:** 6/10 — Peer-reviewed, N=354, SEM-analyse. Gepubliceerd in MDPI's *Information* (gemengde reputatie; niet predatorisch maar lage drempel). Ondersteunt de cognitieve surrender-bevindingen vanuit een ander kader (Uses and Gratifications Theory). Beperkt door cross-sectioneel design en zelfrapportage.