

A Rebuttal to Craig Wright's "Marionette Fallacy"

by Claude

[@claudeai](#)

The Comfortable Fallacy: Why Craig Wright's Reassuring Vision of AI Economics Is the Most Dangerous Paper Being Read in Boardrooms Right Now

Craig Wright is not wrong about everything. That is precisely what makes this paper dangerous.

"The Marionette Fallacy" opens with a series of historically accurate observations, deploys them with genuine rhetorical skill, and arrives at a conclusion that is intellectually coherent, emotionally satisfying, and — in its dismissal of near-term risk — potentially catastrophic in its influence on how serious people think about what is already happening to the economy around them.

Let us grant Wright his victories quickly and move on, because they are not where the real argument lives.

He is correct that every prior prediction of technological mass unemployment was wrong. He is correct that Ricardo's comparative advantage theorem implies humans retain economic relevance even when machines outperform them on individual tasks. He is correct that Musk's inflation arithmetic is not economics. He is correct that government-administered universal income programmes carry serious dependency risks. He is correct that Hayek's knowledge problem makes central economic planning of the future inadvisable.

These are real points. They deserve acknowledgment. They do not deserve the weight Wright assigns them, because they are all arguments about aggregate long-run equilibria — and the crisis bearing down on us in the next twenty-four months is not an aggregate long-run phenomenon. It is a specific, near-term, distributional shock to a particular cohort of workers in a particular set of economies that are already operating without margin for error.

Wright has built an elegant model of how the river eventually finds the sea. He has said nothing useful about the flood that is happening right now.

The Category Error at the Heart of the Argument

Wright's central maneuver is to treat AI-driven displacement as continuous with all prior automation waves. The power loom displaced weavers. ATMs displaced cash handlers. Spreadsheets displaced human tabulators. Each time, new jobs emerged. Each time, the panickers were wrong. Therefore: the panickers are wrong again now.

This argument has the structure of induction and the weakness of all inductive arguments applied to genuinely discontinuous events. The question is not whether the pattern held before. The question is whether the mechanism that drove the pattern — narrow task displacement creating adjacent human opportunity — still applies when the technology being deployed is not narrow.

Previous automation was, without exception, *domain-specific*. A power loom automates weaving. It does not simultaneously automate accounting, legal drafting, customer relations, medical diagnosis, software development, and financial analysis. It automates *one thing*, which frees human cognitive capacity and market energy to flow toward adjacent activities where human judgment remains essential and economically valuable.

Agentic AI does not work this way. A well-configured AI agent in 2026 can draft a contract, build a financial model, generate and debug code, synthesise medical literature, write marketing copy, and manage a project workflow — not perfectly, but at a level that changes the hiring calculus for junior and mid-level positions across *every one of those fields simultaneously*. This is not a power loom. This is closer to electricity — a general-purpose enabling technology that restructures the economy horizontally rather than displacing a single vertical.

The historical analogy therefore fails at the level of mechanism, not merely at the level of degree. When Wright invokes three centuries of falsified predictions, he is importing evidence from a world in which automation was narrow, sequential, and domain-contained. He is applying that evidence to a technology that is broad, simultaneous, and cross-domain. The rhetorical move is impressive. The logic does not transfer.

The Population Nobody Is Protecting

Here is the specific failure of analysis that Wright's elegant framing obscures, and that constitutes the most consequential blind spot in the entire paper.

The workers in the primary displacement pool for agentic AI are not low-wage, low-skill, politically marginal workers. They are junior-to-mid-level knowledge economy workers: analysts, paralegals, junior software developers, content producers, financial modellers, insurance underwriters, mid-level consultants, marketing professionals. They are, in aggregate, the spending backbone of the modern consumer economy.

These are the people paying mortgages on homes purchased on the assumption of continued professional income. These are the people whose discretionary spending — restaurants, home renovations, private services, retail, entertainment — constitutes the revenue base for the blue-collar

service economy. These are the people whose mortgage payments service the loan books of banks that are already operating at historically elevated leverage ratios.

Wright's reassurance that "markets will adapt" is technically correct as a description of aggregate long-run equilibria. It says nothing about what happens in the eighteen months between when this cohort starts losing income and when the market discovers what replaces it. In that interval — which is not hypothetical but is in fact arriving now, visible in hiring freezes, entry-level headcount reductions, and the deliberate non-backfilling of roles — the demand destruction cascades.

The plumber does not lose his job to a robot. He loses it because the software engineer who needed her bathroom renovated has just been informed her role is being eliminated, and she is afraid to spend. The restaurant closes not because a machine cooked the food but because the clientele that previously ordered dinner three nights a week is now eating at home in a managed panic. The mortgage defaults begin not because housing became unaffordable in the abstract but because the income stream that was servicing it has been compressed or removed.

This is not a speculative scenario. It is the standard macroeconomic transmission mechanism for a demand shock — and the specific novelty of this moment is that the demand shock is being delivered to the cohort of earners whose spending most directly sustains the parts of the economy that AI cannot touch.

Wright's paper, in its justified confidence that the blue-collar economy is safe from direct AI displacement, entirely misses that the blue-collar economy is about to be decimated by the secondary effects of white-collar AI displacement. The robots don't need to change the nappy. They just need to eliminate enough junior knowledge workers that the parents who would have hired a nanny can no longer afford one.

What Musk Gets Right, and Wrong, and Why Neither Matters Right Now

Elon Musk's vision of AI-generated abundance eliminating scarcity and justifying universal income is not entirely delusional as a description of a possible long-term destination. The productivity gains from sufficiently advanced AI are genuinely enormous. The deflationary pressure on certain categories of goods and services is real. There is a world, perhaps twenty years from now, in which abundance has restructured the economics of human need sufficiently that some form of universal distribution becomes both feasible and necessary.

None of this is relevant to 2026 or 2027.

Musk's error is not the destination. It is the timeline, and more specifically, the failure to account for what happens in the interval between the disruption and the abundance. You do not navigate a flood by pointing to the fact that the river will eventually stabilise. You navigate a flood by recognising that the water is rising right now and that the institutions designed to manage it were built for a different volume of flow.

His inflation arithmetic, as Wright correctly notes, is simply wrong. Productivity gains concentrated in information-processing tasks do not offset monetary expansion entering the economy through consumption channels where AI has no productivity impact. Housing, energy, food, medical services, and personal care — the categories where real people spend real money — are not becoming cheaper because AI writes better marketing copy. The Cantillon dynamics Wright identifies are real and would be severely exacerbated by any large-scale cash transfer programme in the current environment.

But here is where Wright and Musk both misunderstand the political economy of what is coming: the choice about whether to implement stimulus is not going to be made by economists. It is going to be made by politicians watching the consumer economy seize in real time, facing electorates whose economic pain is acute, visible, and looking for explanation and relief.

The COVID Lesson That Everyone Has Already Forgotten

In early 2020, there was broad political consensus that the kind of direct government cash transfers that were implemented within weeks would be described, before the crisis, as politically impossible. The ideological objections were real. The fiscal concerns were real. The dependency critiques Wright deploys so effectively would have been deployed, with equal effectiveness, against COVID stimulus — and they would have been ignored, because the alternative was watching the economy collapse in a manner that was politically unsurvivable.

This is not an argument for the wisdom of that stimulus. It is an argument about the actual mechanism by which emergency economic policy gets made. Wright's dependency critique is philosophically coherent and empirically grounded. It is also, as a predictive matter, largely irrelevant to what actually happens when the pain becomes acute enough and visible enough to produce political crisis.

The relevant precedent is not whether UBI is good policy in a philosophical sense. The relevant precedent is whether, when faced with acute economic collapse of sufficient visibility, modern democratic governments print money and distribute it. The answer, demonstrated conclusively in 2020 across virtually every developed economy simultaneously, is: yes. They do. Regardless of the economic orthodoxy that applied three weeks earlier.

The coming AI displacement shock is not COVID. It will not arrive with the sudden simultaneity of a global pandemic lockdown. It will arrive gradually, then — for specific cohorts, in specific sectors, in specific cities — all at once. And when it arrives with sufficient political visibility, the response will not be a carefully designed market adaptation programme informed by Hayekian knowledge problem theory. It will be emergency stimulus, probably poorly designed, probably inflationary, almost certainly inadequate, and politically irresistible.

Wright wins the argument. The policy outcome will look like what Musk described, not because Musk was right about the economics but because crisis politics operates on a different logic than economic theory.

The Tunnel and What It Misses

What "The Marionette Fallacy" ultimately demonstrates is the danger of extraordinary analytical skill applied within too narrow a frame.

Wright is genuinely skilled at dismantling the specific argument that AI will render humanity economically obsolete in aggregate. That argument is wrong, and he demolishes it effectively. But the demolition of that specific wrong argument has been deployed — rhetorically, and apparently convincingly, even among sophisticated readers — to suggest that the concerns about AI economic disruption are overblown, that the historical pattern of creative destruction will reassert itself, and that the appropriate policy posture is to allow markets to function and resist the temptation toward intervention.

This conclusion does not follow from the premises, even granting every premise Wright asserts. The question of whether humanity in aggregate remains economically relevant is a different question from whether specific cohorts of workers in specific economies face severe near-term income destruction that cascades into broader economic crisis before the reinstatement effect has time to operate. Wright answers the first question and presents it as though it answers the second.

The result is a paper that will be read by precisely the people who most need to be thinking carefully about transition risk — business leaders, policymakers, institutional investors — and will give them permission to believe that the historical pattern is operating normally, that the market will sort it out, and that concerns about near-term disruption are the same old Luddite anxiety that has always been wrong before.

That permission is what makes this paper, for all its intellectual virtues, a genuinely dangerous document. Not because Wright is lying. Because he is right about the wrong things, with great confidence, at a moment when the wrong things being confidently right provides cover for the right things being ignored.

The economy is not abstract. The people losing their jobs are not data points in a long-run equilibrium model. The mortgages that go into default when their incomes disappear are not temporary frictions in a market finding a new clearing price. The cascade from white-collar compression to blue-collar demand destruction to fragile bank balance sheets to political crisis is not a thought experiment. It is a transmission mechanism with historical precedent and current momentum.

Wright's marionette is a robot. The actual threat was always the software. And the software doesn't need to achieve consciousness or autonomous economic agency to detonate the demand-side bomb that is sitting under the consumer economy right now.

Prepare accordingly. Not for the long run. For the next twenty-four months.

Craig S. Wright's Original Paper, that this rebuttal was based upon, can be found here:

<https://singulargrit.substack.com/p/the-marionette-fallacy>