Why neither humans nor Al alone match the performance of a deliberately paired system – and how the human provides the vector and guardrails within which Al excels.

October 2025



About the Canadian Centre for Economic Analysis

The Canadian Centre for Economic Analysis (CANCEA) is a premier socio-economic research and data firm, distinguished for its unwavering commitment to delivering objective, evidence-based analysis. Anchored in a holistic understanding of market shifts, policy implications, and economic behaviors, CANCEA's research transcends traditional boundaries to offer a panoramic view of the socio-economic landscape.

Driven by modern data science techniques, including the pioneering use of agent-based modelling, CANCEA's analytical spectrum encompasses a diverse range of services. The cornerstone of CANCEA's analytical prowess is its state-of-the-art agent-based platform, the largest in North America, that is a meticulously crafted data-driven model that encapsulates over 56,000 distinct regions across Canada. This platform facilitates in-depth, multidisciplinary analyses, empowering stakeholders with unparalleled insights into the interplay of various socio-economic parameters.

Embracing a systems-centric approach, CANCEA uniquely adopts a single-model strategy. This integrated approach allows for the seamless fusion of multiple disciplines and stakeholder perspectives, culminating in holistic, collaborative, and quantitative analyses that inform and guide pivotal market, policy, and economic decisions.

© 2025 Canadian Centre for Economic Analysis

Printed in Canada • All rights reserved

Moreover, as champions of data integrity and comprehensiveness, CANCEA offers robust Canadian data services, ensuring that stakeholders are equipped with the most accurate and up-to-date information for informed decision-making.

In essence, CANCEA is at the forefront of socio-economic research, transforming complex data into actionable insights for a diverse range of sectors and stakeholders.

About the Author

Paul Smetanin is Founder, President and CEO of CANCEA. A 35-year veteran of systems modelling, risk management and complexsystems economics, he leads Canada's integrated largest, socio-economic simulation platform to produce decisiongrade evidence for governments and industry. His work spans infrastructure, housing, labour markets, healthcare, impact analysis and more, with over publications. Paul has served as a C-suite executive at ANZ and Algorithmics, before building CANCEA's ONEMODEL™ paradigm and its explainable 5W Attribution™ method. These capabilities are productized as an enterprise SaaS that bridges agent-based simulation with modern AI/LLM workflows advancing trustworthy, transparent, machine-assisted policy and strategy.

Citation: Smetanin, P. The Conductor Model: Direction with Magnitude for Human-Al Teams. Canadian Centre for Economic Analysis. October 2025.



ABSTRACT

Modern AI multiplies magnitude - the volume and speed of remedial and iterative work - yet it neither chooses ends nor reconciles values. This paper advances the Conductor Model, in which accountable humans supply direction - purpose, values, constraints, trade-offs, measurable objectives, and rigour - while AI provides scalable execution within those rails. We explain why "human versus AI" is the wrong question; map what is non-delegable to humans; and show how human oversight protects causality in systems that are like "supercharged autocomplete" unless told otherwise. We summarise what AI is - and isn't - good at; the risks of ignoring the Conductor role (automation bias, manual pride, provenance drift, metric gaming); and why conceptual knowledge serves as the loss function that penalises wrong answers. We introduce, in brief only, the V.E.C.T.O.R. approach (Vision, Ethics, Constraints, Trade-offs, Objectives, Rigour) as an operationalisation of the Conductor Model, with full templates and checklists reserved for a companion paper. The claim is not romanticism about "augmentation"; it is an operational stance about throughput, quality control, and accountability.

Audience & scope: This paper is written for senior executives, programme leads, and technical leads who must decide who does what in human—AI work. It is an operational playbook, not a methods tutorial; formal templates, proofs, and implementation checklists live in the companion paper.

Direction without magnitude stalls; magnitude without direction wanders. Together they form a usable vector.

About the author: Paul Smetanin, Founder and CEO of CANCEA, is a key opinion leader on complex systems analysis, applying machine-learning—informed, agent-based models to societal, market, policy planning and stress testing and risk management. For 25 years he has advanced ONEMODEL (4,000+topics), with 5W Attribution revealing the "why." He is productising as SaaS and, drawing on a 35-year career, speaks on construction and infrastructure governance.



TABLE OF CONTENTS

| Abstra | ct | 2 |
|--------|--|----|
| 1.0 | Introduction | 4 |
| 1.1 | Definitions Used In This Paper | 5 |
| 1.2 | How ABM, LLMs, and ML Fit Together in the Conductor Model | 6 |
| 2.0 | Why "Human vs. AI" Is the Wrong Question | 7 |
| 3.0 | The Conductor Role & What Is Non-delegable | |
| 3.1 | Non-Delegable Responsibilities (human, by design) | |
| 3.2 | Delegable work: A Simple Heuristic | 8 |
| 4.0 | Protecting Causality Under Human Oversight | 10 |
| 4.1 | Make Causality Explicit - Label Claims | |
| 4.2 | Set Identification and Calibration Standards Up Front | |
| 4.3 | Preserve Provenance by Construction ("Source-Lock") | |
| 4.4 | Demand Reproducibility and Uncertainty | 11 |
| 4.5 | Keep a Decision Log and Change Control | 11 |
| 4.6 | Minimal Guardrails that Specifically Protect Causality | 11 |
| 5.0 | What Al Is (and Isn't) Good At | 12 |
| 5.1 | Strengths: Scalable Magnitude on Well Specified Work | 12 |
| 5.2 | Limits: Direction, Trade-Offs, and Caution Cannot Be Delegated | 13 |
| 6.0 | The Risks of Ignoring the Conductor Model | 14 |
| 7.0 | Why Conceptual Knowledge Is the "Loss Function" | 15 |
| 7.1 | Conceptual Knowledge, Operationalised | 15 |
| 7.2 | How Errors Amplify Without Concept Mastery | 15 |
| 7.3 | Embedding Conceptual Knowledge so the Machine Stays Inside the Rails | 15 |
| 8.0 | Preview: The Vector Approach as the Authorisation | 17 |
| 8.1 | The Core Artefact: A Short Vector Brief (2–3 pages) | 17 |
| 8.2 | Minimal Workflow (Direction First, Then Magnitude) | 17 |
| 8.3 | Guardrails Baked Into V.E.C.T.O.R. | |
| 8.4 | Authority and Accountability (Who Signs What) | |
| 8.5 | Stop/Go Triggers (Summary) | 18 |
| 9.0 | Worked example — Housing delivery: ABM + LLM/ML Under V.E.C.T.O.R | |
| 9.1 | Why This Example Appears Here | |
| 9.2 | The Setup | |
| 9.3 | Two Compounding Problems | |
| 9.4 | What the Top-Down View Ignores: The Housing Pipeline | |
| 9.5 | Bottom-Up Alternative: Agent-Based Modelling (ABM) With Guardrails | |
| 9.6 | Direction First: A Two-Page Vector Brief (expanded excerpt) | |
| 9.7 | Magnitude Next: What Al/ML Can Do Safely | |
| 9.8 | Outcome (Why the Numbers Change) | |
| 9.9 | Why This Supports the Paper's Argument | |
| 10.0 | Conclusion: Leadership with Leverage | |
| A. B | ibliography | 25 |



1.0 INTRODUCTION

The question most executives and team leads still ask about artificial intelligence is framed as a contest: Al or human - who wins? It is the wrong question. Framed that way, organisations either abdicate judgement to tools that do not possess it, or they block tools that could multiply the value of their people. The better question is: What is the right division of cognitive labour between humans and Al so that together they produce outcomes neither can achieve alone? This paper argues for a human-in-command stance in which concept-literate practitioners act as conductors - setting objectives, constraints, and standards - while Al provides scalable magnitude to carry that intent across the many remedial and iterative tasks that turn vision into verified outputs. Direction without magnitude stalls; magnitude without direction wanders. Together they form a usable vector.

Al has changed daily knowledge work - not by replacing decisions, but by compressing the time and cost of everything between an idea and a decision-ready artefact. Literature sweeps that once took days now take hours; formatting and assembly that previously needed multiple passes can be completed in minutes. None of that is useful unless someone first defines what counts as good and why it matters. In our framing, the thesis is simple: Humans supply direction; Al supplies scalable magnitude. When we behave as conductors - explicitly setting purpose, values, constraints, trade-offs, measurable objectives, and rigour tests - Al can handle the high-volume work that makes complex projects feasible within real world budgets and timelines. This is not a romantic claim about "augmented intelligence"; it is an operational claim about throughput, quality control, and accountability.

Two pathologies make the human versus AI framing especially costly. Automation bias is the quiet temptation to outsource goal selection and value judgements to systems optimised to continue patterns, not to adjudicate among purposes. Manual pride is a professional identity organised around artisan workflows that cannot scale to the scope and speed demanded of contemporary work. Both create mismatches between the magnitude of the task and the resources available. The Conductor Model resolves this tension by separating the choosing of ends from the executing of means - and making that separation explicit, documented, and auditable. (Parasuraman & Riley, 1997; Skitka, Mosier, & Burdick, 1999)

This paper has three pragmatic aims. First, name the strengths and limits of current AI systems in language operational leaders can use - what kinds of work can be delegated with confidence, what cannot, and what must be done jointly. Second, offer a Conductor's framework for direction - V.E.C.T.O.R. (Vision, Ethics, Constraints, Trade-offs, Objectives, Rigour) - as both a mental model and a checklist that can be embedded in briefs and decision logs. Third, demonstrate the division of labour in domain neutral "workflow moments," showing where AI reduces cycle time and cognitive load, and where human judgement is non-delegable, especially in protecting causality in a probabilistic, correlation driven technology. In this paper we introduce V.E.C.T.O.R. only in summary as the operationalisation of the Conductor Model; full templates, acceptance test gates, and role separated controls appear in a companion paper.



The Conductor stance requires conceptual literacy, not expertise in gradient descent or transformer internals. It is the mastery of the concepts that govern the work itself - mechanisms, admissible measures, and limits - that allows leaders to specify direction and to evaluate whether AI produced artefacts are fit for purpose. Absent those lenses, AI accelerates confusion; present them, and AI accelerates execution. Equally, the stance demands explicit standards: when AI enters the loop, tacit norms must be written down so they can be embedded in prompts, acceptance criteria, and tests. This is not bureaucratic overhead; it is the price of reliability at scale.

By the end of the paper, a chief executive, programme lead, or technical lead should be able to (a) recognise where AI belongs in a workflow; (b) specify direction using V.E.C.T.O.R.; (c) institute guardrails that reduce automation bias and preserve provenance; and (d) measure quality in ways that are reproducible and auditable. This is not about adopting a tool; it is about adopting a role. The conductor is not a bottleneck; the conductor creates usable vectors - pairing direction with magnitude - so that organisations can move not merely faster, but on purpose.

Who this paper is for (and how to read it)

- Primary: decision-makers (executives, programme leads) who need a credible division of labour: humans supply direction; Al supplies magnitude.
- Secondary: technical leads who must operationalise guardrails and acceptance tests.
- What this paper is not: a methods or algorithms tutorial. The detailed templates and checklists for V.E.C.T.O.R. are in a companion paper.
- Reading path: Sections 2–8 are decision-focused; technical readers can jump to the call-outs labelled "For technical readers," and to the worked example (Section 9).

1.1 DEFINITIONS USED IN THIS PAPER

Human (Conductor): The accountable agent who sets purpose, values, constraints, and quality standards; makes trade-offs; and bears responsibility for decisions and outcomes.

Al (Model/System): A general-purpose probabilistic system that performs pattern completion and transformation across representations (text, code, tables, images) without intrinsic goals or values.

Remedial Tasks: High-volume, iterative activities that apply established criteria to known inputs (e.g., literature extraction, summarisation, formatting, boilerplate assembly, code scaffolding, data-cleaning stubs).

Direction/Vector: Direction comprises purpose, values, constraints, and acceptance criteria set by the human; the vector is the combination of that direction with the magnitude AI supplies through scalable execution.



Machine learning (ML): A set of methods that learn patterns from data to make predictions, classifications, or estimates by optimising an objective on examples rather than by hand-coding rules. In practice, ML produces a function that maps inputs to outputs and is evaluated by how well it generalises to new data.

Large language models (LLMs): Very large neural networks trained on massive text corpora to predict the next token in sequence. Because of this training, they can generate and transform text, follow instructions, and structure information across many tasks. Their outputs are probabilistic and reflect statistical patterns in the training data; they require external constraints and checks where truth, causality, or safety matter.

1.2 How ABM, LLMs, and ML Fit Together in the Conductor Model

We use agent-based modelling (ABM) as the running testbed because it encodes mechanisms and stock-flow constraints by construction: agents, queues, and capacity gates make conservation and lags auditable. In our division of labour, LLMs and ML supply magnitude—rapid extraction, formatting, code scaffolding, calibration fits, and scenario enumeration—under the rails set by the human conductor and the model's structure. This pairing protects causality and provenance in a technology that otherwise optimises for plausibility. (See Sections 4 and 7.)

To make this concrete, think of the work as a simple relay: the human sets the destination and the safety rules; ABM lays down the rails; ML tunes the knobs; LLMs move the paperwork and scaffolding quickly; and the human signs off. Here's how those hand-offs feel in practice:

- First, set the destination (human). We decide the purpose, audience, and the decision we're informing. We also write the acceptance tests up front—what evidence will count as "credible" when we're done.
- **Next, lay the rails (ABM).** We build the mechanism: agents, states, queues, capacities, and timing. This forces stock-and-flow conservation and makes each state change auditable, so only causally admissible moves are even possible in the sandbox.
- Then, tune the knobs (ML). Using data, we estimate the parameters the mechanism needs—things like approval lead times, failure rates, or demand shocks—and attach uncertainty bands, so scenarios have defensible ranges rather than guesses.
- Meanwhile, move the paperwork fast (LLMs). We use LLMs to extract and normalise text, assemble documents to a schema, scaffold code, and spin up scenario variants—high throughput work that stays inside the rails we've already set.
- **Finally, gate the quality (human).** We check results against the acceptance tests, label causal vs. associative claims, confirm provenance and reproducibility, and sign the decision log.



2.0 WHY "HUMAN VS. AI" IS THE WRONG QUESTION

Purpose: reframe the problem from a contest to a division of labour and justify a human-in-command stance.

The "human vs. AI—who wins?" framing is the wrong one. When treated as a contest, organisations either hand judgment to tools that don't have it, or they shut out tools that could amplify their people's impact. The right question is how to split the cognitive work so humans and AI achieve results neither could deliver alone. In our view, accountable humans set direction—purpose, values, constraints, trade-offs, measurable objectives, and rigour—while AI provides scalable magnitude on the remedial and iterative tasks that carry intent through to execution. Direction without magnitude stalls; magnitude without direction drifts. Put together, they create a usable vector.

This reframing is practical, not rhetorical. All has changed daily knowledge work by compressing the time and cost of everything between an idea and a decision ready artefact - from literature sweeps and structured extraction to formatting and code scaffolding. Yet none of that speed is useful unless someone first defines what counts as good and why it matters. The Conductor Model therefore separates the choosing of ends from the executing of means, making that separation explicit, documented, and auditable.

Two recurring pathologies explain why the "contest framing" is costly. Automation bias tempts teams to outsource goal selection and value judgements to systems optimised to continue patterns, not to adjudicate among purposes. Manual pride clings to artisan workflows that cannot scale to contemporary scope and speed. Both create a mismatch between the magnitude of the task and available resources. The conductor stance resolves this tension by pairing human direction with machine magnitude under explicit standards, provenance, and tests.



3.0 THE CONDUCTOR ROLE & WHAT IS NON-DELEGABLE

Purpose: define non-delegable, delegable, and joint responsibilities, and provide a practical heuristic.

The Conductor Model is a work discipline: humans set direction; Al delivers magnitude. The conductor is the accountable agent who specifies purpose, values, constraints, trade-offs, objectives, and rigour - and then delegates high volume, iterative tasks to Al within those rails. Done well, throughput rises without diluting judgement or accountability.

3.1 Non-Delegable Responsibilities (Human, By Design)

- Purpose & audience (Vision): Name the decision, the decision-maker, and the decision date; articulate the minimally sufficient artefact and the thesis it must defend.
- Values & duty of care (Ethics): Make distributional and fairness lenses explicit; bound harms; decide which claims require higher bars of evidence.
- Boundaries (Constraints): Fix time, data rights, and compliance requirements including allowed sources and accessibility obligations.
- Knob settings (Trade-offs): Choose breadth vs. depth and speed vs. precision; record the rationale that makes these choices fit for purpose.
- Acceptance bars (Objectives): Convert intent into tests quantitative and qualitative specifying who tests, when, and how.
- Accountability (Rigour): Maintain decision logs, adversarial checks, provenance tables, and reproducibility so changes and claims can be explained and rerun.

Causality remains a human responsibility. Modern AI is a probabilistic, correlation driven pattern recogniser. Causal claims require human oversight - with identification strategies (or documented calibration), explicit uncertainty, and honest labelling of causal/associative/speculative.

3.2 Delegable work: A Simple Heuristic

Delegable work (to AI, under rails) include remedial and iterative tasks - structured extraction, summarisation to rubric, formatting and document assembly, boilerplate drafting, code scaffolding, combinatorial exploration, and coordinated summarisation - are well-suited to AI once the rails are set. The value is not only speed but consistency against a defined schema.

Joint work (human + AI) is where the task decomposes into machine doable steps but requires human synthesis at the end - e.g., scenario enumeration followed by selection, or reteam prompt generation followed by adjudication - the most effective pattern is AI assisted breadth \rightarrow human choice and narrative.

When direction (human) meets magnitude (AI) under explicit standards and provenance, teams produce usable vectors - work that moves faster and on purpose, with accountability intact.

A simple heuristics is to delegate to AI when you can specify the input format, the allowed sources, and the acceptance tests, and when failure has low externalities and is cheap to detect. Retain human control



when tasks define what counts as good or require reconciling values, risks, and consequences. Use joint workflows whenever work can be decomposed into machine doable steps with human set standards and human synthesis at the end.



4.0 PROTECTING CAUSALITY UNDER HUMAN OVERSIGHT

Purpose: explain why causality is a non-delegable human responsibility and how to set acceptable evidence standards in practice.

Modern AI is a probabilistic, correlation driven pattern recogniser. Left unattended, it will produce fluent, plausible continuations of patterns - not warranted causal claims. Protecting causality is therefore a non-delegable human responsibility: identify where causal language is permitted, demand appropriate evidence or calibration, and label claims honestly. (Pearl, 2009; Hernán & Robins, 2020).

LLMs pose a distinctive risk: as conditional, probabilistic pattern-generators, they optimise for plausibility, not proof. Their fluency plays to our brain's love of cohesive narratives and rhythmic phrasing, so they can sound deeply convincing while echoing correlations rather than verified cause-and-effect. It's like a song with an irresistible beat and a sing-along chorus, even if the lyrics don't make sense. The tune lands: the truth doesn't, so treat their output as catchy drafts that need checking when accuracy matters.

Direction supplies the rules of inference; *magnitude* executes within those rules.

4.1 Make Causality Explicit - Label Claims

Adopt a simple labelling discipline for every material statement:

- Causal: allowed only when an identification or calibrated design is present (e.g., experiment, quasi experiment, structural model with sensitivity analyses).
- Associative: relationships supported by patterns and controls, without identification.
- Speculative: hypothesis or mechanism sketch, marked as such.

These labels travel with the claim through drafting, review, and publication so that readers know what kind of warrant they are being asked to accept.

4.2 SET IDENTIFICATION AND CALIBRATION STANDARDS UP FRONT

Write, in plain language, the admissible methods for causal claims in this work: what qualifies as identification, when calibration is acceptable, and which sensitivity checks must pass. Treat this as part of the brief (the Rigour and Objectives in V.E.C.T.O.R.), not an afterthought. Where causal warrant is out of scope, require the system to stay associative and label it accordingly. (Hernán & Robins, 2020; Pearl, 2009)

4.3 Preserve Provenance by Construction ("Source-Lock")

Require models to cite from a permitted source pack and to emit a provenance table mapping each material claim to its supporting source or reproducible code. Expansion of sources is a recorded decision, not an accident. Provenance turns speed into auditability. (W3C, 2013; Office for Statistics Regulation, 2021)



4.4 DEMAND REPRODUCIBILITY AND UNCERTAINTY

Causal and quantitative artefacts must be rerunnable from a clean environment (seed control, environment spec) and must state uncertainty (ranges or qualifiers), not just points. Qualitative claims include confidence qualifiers tied to evidence strength. Reproducibility and uncertainty statements are part of the acceptance tests at the quality gate¹. (Office for Statistics Regulation, 2021)

4.5 KEEP A DECISION LOG AND CHANGE CONTROL

Maintain a decision log that records what changed, why, and on whose authority - especially changes to Vision, Ethics, Constraints, Trade-offs, Objectives, or Rigour. Tie each change to the relevant V.E.C.T.O.R. element so reviewers can see how direction evolved and whether causal warrants were upgraded, downgraded, or deferred.

4.6 MINIMAL GUARDRAILS THAT SPECIFICALLY PROTECT CAUSALITY

Enforce role separation: Author ≠ Checker ≠ Approver.

The conductor is responsible for direction and decisions; a separate reviewer is responsible for provenance and data rights compliance. This preserves accountability while leveraging AI for magnitude within explicit rails. The minimal guardrails that specifically protect causality include:

- Automation bias checks: ensure no step asks AI to choose goals or make causal leaps beyond the allowed warrant; verify evidence to claim mapping.
- Adversarial prompts: ask for counter theses, minimal change failures, and Goodhart style failure modes; require citations from the source pack.
- Prompt-injection/data-exfiltration defenses. Treat untrusted inputs (docs, URLs, user content) as
 adversarial; restrict tool-use and retrieval scopes; and apply the "OWASP Top 10 for LLM
 Applications" to your acceptance tests (OWASP Foundation, 2025).
- Red team questions: identify which single claim, if false, would most harm credibility; what evidence would reverse the recommendation; and where the pipeline's single point of failure sits.

¹ A quality gate is a deliberate checkpoint in a workflow where work must meet specific, pre-agreed standards before it's allowed to move forward.



5.0 WHAT AI IS (AND ISN'T) GOOD AT

Purpose: catalogue strengths and limits to guide safe delegation and joint work.

Leaders do not need romance about "intelligence"; they need operational clarity: what to delegate, what to retain, and what to do jointly. A practical rule holds across domains: Al excels when acceptance criteria are objective and testable; it struggles when tasks require selecting goals, reconciling values, or exercising epistemic caution under uncertainty.

5.1 STRENGTHS: SCALABLE MAGNITUDE ON WELL SPECIFIED WORK

1. Recall at scale and structured extraction.

Models can sweep large corpora to assemble candidate facts, definitions, and citations, then fit them into consistent structures (literature grids, requirement matrices, data dictionaries). The value is not just speed; it is consistency once a schema is defined.

2. Pattern completion across formats.

Given exemplars, AI produces "more of the same": convert notes to bullet summaries; normalise units and categories; map fields between schemas; draft figure captions from data descriptions; generate test cases from acceptance criteria.

3. Multitype drafting and redrafting.

All efficiently produces variants to a brief - different audiences, lengths, tones - and outline scaffolds that turn a blank page into sections with thesis sentences.

4. Formatting and document assembly.

From reference lists to figure call outs and house-style headers to alt text, models are increasing there potential to enforce deterministic rules and generate boilerplate, reducing error and freeing human attention for argument and design.

5. Code generation and data prep scaffolding.

For "glue code" and routine analytics, AI is a force multiplier: cleaning scripts, reproducible notebooks, plotting functions with docstrings, unit tests, and CI snippets - plus validators and batch run stubs aligned to your schema.

6. Combinatorial exploration.

When breadth is needed - enumerating scenarios, options, risks - AI supplies it quickly, generating alternatives that humans filter and recombine.

7. Co-ordinated summarisation.

Given an explicit rubric, AI compresses content faithfully (e.g., "200 word synopsis against this rubric"), provided outputs are spot checked.



5.2 Limits: Direction, Trade-Offs, and Caution Cannot Be Delegated

1. Goal selection and problem framing.

Models do not choose ends. Asking AI to propose objectives' risks adopting implicit values embedded in training data.

2. Cross domain trade offs.

Balancing multiple criteria and long-horizon effects is judgement laden and context specific. Al can list trade offs; it cannot adjudicate them for you.

3. Value judgements and distributional consequences.

Fairness, duty of care, and acceptable harm require explicit ethical reasoning and accountability, not pattern matching.

4. Epistemic caution and calibration.

Al are a fluent guesser: without tight constraints, it may state plausible falsehoods, overgeneralise, blur provenance, and fail to know when it does not know - dangerous where error bars matter.

5. Theory, causality, and identification.

Summarising literature is not understanding mechanisms. Al does not reliably distinguish correlation from causation or select identification strategies. Credible causal claims require domain theory, defensible assumptions, and tests.

6. Out-of-distribution edge cases.

With sparse, noisy, or novel data, models revert to priors - precisely where human expertise must lead. (Koh et al., 2021; Hendrycks & Gimpel, 2017; Hendrycks et al., 2019)

7. Jagged frontier

Even within the same project, some tasks lie firmly inside current model strengths while adjacent tasks do not. Outside that frontier, human performance can *degrade* when aided by Al. Conductor-set acceptance tests and role separation are the antidote (Dell'Acqua et al., 2023)

8. Long-horizon coherence and accountability.

Sustaining a single, coherent argument across time - and owning the decision - is a human responsibility. Al can help keep a decision log; it cannot sign it.

9. Legal, data rights, and confidentiality constraints.

Models cannot guarantee compliance or obtain permissions. Humans must set and enforce boundaries within which AI operates.

Use the simple heuristic set out earlier: delegate when inputs, sources, and acceptance tests are explicit and failure is cheap to detect; retain when tasks define what counts as good or reconcile values see the conomic analysis. The end. This keeps Al squarely in the magnitude seat and humans in command of direction.

6.0 THE RISKS OF IGNORING THE CONDUCTOR MODEL

Purpose: name recurring failure modes without the conductor role and the structural remedies.

When the conductor role is absent, speed multiplies error. Organisations either abdicate judgement to tools that do not possess it (automation bias) or cling to artisan workflows that cannot scale (manual pride). Both failure modes create a mismatch between the magnitude of the task and the resources available. The Conductor Model resolves this by separating the choosing of ends from the executing of means - and making that separation explicit, documented, and auditable. Without it, the system drifts.

Six recurring risks follow when direction is not fixed before magnitude is applied:

- 1. Accountability drift. Decisions are justified by "what the model produced" rather than by a human rationale tied to purpose, values, and constraints; author, checker, and approver collapse into a single undifferentiated role. Remedy: keep role separation and a signed decision log as non negotiables.
- 2. Provenance slippage. Claims float free of permitted sources, drafts commingle grey literature with primary evidence, and figures cannot be regenerated on demand. Remedy: require "source lock" to a curated source pack and provenance by construction (every material statement mapped to a source or reproducible code).
- 3. Metric gaming (Goodhart). Optimising an easy proxy degrades the underlying objective; eloquence displaces evidence. Remedy: have the conductor define acceptance bars ex ante and run adversarial checks ("show how optimising our metric could harm the objective; propose safeguards"). (Goodhart, 1975; Strathern, 1997)
- 4. Causality overstatement. Correlation driven outputs are presented as causal claims, with no identification or calibration and no uncertainty statement. Remedy: enforce causality labelling (Causal/Associative/Speculative) and admit causal language only under specified warrants and sensitivity tests.
- 5. Scope creep and normative slippage. Systems introduce sources, domains, or value frames outside the permitted envelope, smuggling in unstated winners and losers. Remedy: declare Ethics and Constraints up front (V.E.C.T.O.R.), and quarantine any out of envelope content pending an explicit, logged decision to expand scope.
- 6. Irreproducible throughput. Teams move faster but cannot rerun, defend, or evolve the work; one person's "prompt magic" becomes a single point of failure. Remedy: make reproducibility a quality gate (seed control, environment spec, rerun instructions) and treat tacit standards as explicit rigour requirements.

Early warning indicators of these risks include missing acceptance tests, unlabelled causal claims, outputs citing outside the source pack, explanations that lean on model eloquence rather than evidence, and absent or stale decision logs. The Conductor Model - codified with V.E.C.T.O.R. as a light touch brief - prevents these pathologies by fixing direction before machines generate magnitude and by keeping human oversight where it belongs.



7.0 WHY CONCEPTUAL KNOWLEDGE IS THE "LOSS FUNCTION"

Purpose: show how conceptual knowledge constrains outputs and becomes the practical loss function.

Tools are fast; concepts make them useful. In this framing, concept mastery functions as the loss function for human—AI work: it constrains what the system is allowed to say and penalises outputs that violate mechanisms, admissible measures, or known limits. In practice, the penalty is concrete - an artefact fails an acceptance test and must be reworked. Absent this discipline, AI's eloquence becomes a force multiplier on error; present it, and the same system becomes a disciplined instrument.

7.1 CONCEPTUAL KNOWLEDGE, OPERATIONALISED

- Mechanisms: the causal stories that explain how levers affect outcomes.
- Measures: what counts as evidence and how it is compared.
- Limits: invariants, legal bounds, and standards of proof that cannot be crossed.

Turning these into acceptance tests (e.g., unit checks, direction of effect checks, distributional coverage, evidence to claim mapping) makes wrong answers detectable by design rather than by luck.

7.2 How Errors Amplify Without Concept Mastery

- Category errors. Stocks are confused with flows; outcomes with controls; reliability collapsed into averages.
- Proxy traps. Convenient metrics stand in for welfare relevant objectives, proxy improvement masquerades as real progress.
- Path dependence. Priors from elsewhere are imported without boundary conditions.
- Goodhart effects. Optimising the metric degrades the objective.
- Normative slippage. "Neutral" framings embed unexamined value choices about winners and losers.

These are not the model "being wrong" so much as the human failing to supply the loss function - the theory and method that render certain outputs unacceptable.

7.3 EMBEDDING CONCEPTUAL KNOWLEDGE SO THE MACHINE STAYS INSIDE THE RAILS

ABM is a practical way to embed these invariants and admissible measures so that LLM/ML magnitude operates inside a conserved state space.

- 1. Concept map first. Name outcomes, drivers, levers, constraints, and observables; annotate edges with expected sign and confidence.
- 2. Define invariants. Conservation identities, budget constraints, legal limits, and unit balances that must hold.
- 3. State admissible measures. Primary metrics, permitted proxies, and where proxies are forbidden.



- 4. Method rules. When causal language is allowed; identification or calibration strategies; strata required for distributional reporting.
- 5. Retrieval-augmented generation (RAG). Where models must ground claims in a permitted *source pack*, route generation through RAG so drafts cite retrieved passages from approved corpora at generation time. This turns 'source-lock' from policy into mechanism (Lewis et al, 2020).
- 6. Acceptance tests. Reproducibility, source traceability, uncertainty statements, and explicit "don'ts" tied to known failure modes.

Once these are explicit, Al's magnitude compounds value rather than compounding error: the system produces more, faster, inside a conceptual envelope that the conductor has set and will defend. This is the practical reason conceptual literacy - not model internals - is non-delegable in high functioning human—Al teams.



8.0 PREVIEW: THE VECTOR APPROACH AS THE AUTHORISATION

Purpose: introduce V.E.C.T.O.R. as the operational authorisation that fixes direction before magnitude.

The Vector Approach (V.E.C.T.O.R.) is the authorisation and operating envelope for human—AI work. It converts the Conductor Model into a light, pre commitment compact: the human conductor sets Vision, Ethics, Constraints, Trade offs, Objectives, and Rigour before any substantive AI runs, and those same elements become the tests applied after the runs. In short, V.E.C.T.O.R. fixes direction so that AI can deliver magnitude without eroding accountability. This section introduces the concept and its guardrails only; full templates, checklists, and role patterns appear in the companion paper.

8.1 THE CORE ARTEFACT: A SHORT VECTOR BRIEF (2-3 PAGES)

The Vector Brief is the authorisation document. It names the decision, decisionmaker, and date (Vision); states value lenses and duty-free boundaries (Ethics); sets permitted sources, data rights, time, and style constraints (Constraints); records intentional knob settings (Trade-offs); translates intent into measurable acceptance bars (Objectives); and defines the accountability machinery - provenance, tests, and decision logging (Rigour). Nothing material runs until the conductor signs this brief. It is the difference between "ask the model to help" and "authorise the model to operate within rails." When the work involves simulation (e.g., ABM), the Brief must state the invariants and conservation tests that outputs must satisfy (e.g., stock-and-flow identities), so that Al magnitude runs inside a conserved state space.

8.2 MINIMAL WORKFLOW (DIRECTION FIRST, THEN MAGNITUDE)

- 1. Vector Brief. Conductor sets V.E.C.T.O.R as above.
- 2. Source Pack. Curate right cleared, permitted sources and the house style.
- 3. Decomposition. Split machine doable tasks (extraction, formatting, scaffolding) from human only tasks (framing, adjudication, synthesis).
- 4. Al Runs (under tests). Execute remedial/iterative work with acceptance checks embedded.
- 5. Human Synthesis. Recompose outputs into the argument; make the trade offs; write the rationale.
- 6. Quality Gate. Apply tests; if a bar fails, adjust Trade-offs or Objectives explicitly and rerun.

This is a thin spine - just enough process to keep speed from outpacing judgement.

8.3 GUARDRAILS BAKED INTO V.E.C.T.O.R.

Source & provenance by construction. Models cite only from the permitted Source Pack; every
material claim maps to a source or reproducible code; expansions of scope are logged. This turns
velocity into auditability. (NIST, 2023; ISO/IEC, 2023; OECD, 2024; European Union, 2024)



- Causality labelling. Claims are tagged Causal/Associative/Speculative; causal language is admissible only under pre stated warrants (identification or documented calibration) with uncertainty stated.
- Acceptance tests as gates. Objectives are measurable; gates fail closed (no publish/ship if tests fail). The model can propose, but only humans pass gates.
- Role separation. Author ≠ Checker ≠ Approver. The conductor signs direction and the decision log;
 a separate reviewer signs provenance/data rights; approver authorises release. All assists but cannot "own" a decision.
- Uncertainty & reproducibility. Outputs carry ranges/qualifiers and can be rerun from a clean environment (seed control, environment spec, rerun instructions).
- Metric discipline (Goodhart guard). Trade offs name the true objective; tests include an explicit "show how optimising our metric could harm the objective" check.
- Data rights & confidentiality boundaries. Constraints enumerate what is permitted; violations are treated as scope breaches requiring explicit re-authorisation.
- Infrastructure & energy constraints. Explicitly enumerate compute quotas, latency/SLA requirements, model-hosting location, and power availability as constraints; these have become first-order determinants of feasibility as data-center electricity demand climbs sharply (IEA, 2024).

8.4 AUTHORITY AND ACCOUNTABILITY (WHO SIGNS WHAT)

- Conductor (accountable human). Signs the Vector Brief and quality gate decision; owns Vision/Ethics/Trade offs and the rationale.
- Checker (independent reviewer). Signs provenance and data rights compliance; verifies causality labels and uncertainty statements.
- Approver. Authorises release to stakeholders; confirms that Objectives/Rigour have passed.
- Al system. Executes within rails; may draft, collate, and test; cannot set ends, relax constraints, or sign.

This preserves human responsibility for ends while exploiting machine magnitude for means.

8.5 Stop/Go Triggers (SUMMARY)

- If any V.E.C.T.O.R element is blank → Stop; complete the brief.
- If outputs cite outside the Source Pack or lack provenance/uncertainty tags → Stop; quarantine and log any scope change.
- If a claim labelled Causal lacks identification/calibration and sensitivity checks → Relabel or remove before proceeding.

These triggers keep the system inside its authorised operating envelope with minimal ceremony.

What follows (in the companion paper) is the full V.E.C.T.O.R. checklist, role separated workflows, acceptance test catalogues, adversarial prompts, and reproducibility runbooks. Here, our purpose is



narrower: establish that V.E.C.T.O.R. is the authorisation mechanism that fixes direction, guards causality and provenance, and lets AI deliver magnitude on purpose.



9.0 WORKED EXAMPLE — HOUSING DELIVERY: ABM + LLM/ML UNDER V.E.C.T.O.R.

9.1 WHY THIS EXAMPLE APPEARS HERE

Up to now we have argued for human direction with machine magnitude, and for causality/provenance guardrails. We now instantiate that stance: ABM encodes the pipeline physics (directional rails); LLM/ML provide magnitude (extraction, calibration, assembly) within those rails.

9.2 THE SETUP

A city announces a policy meant to "unlock" housing - say, a purchase credit or a fee holiday. A top down, aggregate model (sometimes dressed up with machine learning fits to historical averages) projects +10,000 homes in two years: cheaper to buy \Rightarrow demand rises \Rightarrow builders respond \Rightarrow equilibrium restores. The projection looks clean but never shows where those homes come from, who gets them, or what must be true in the pipeline to deliver them. It also does not enforce stickhandle conservation: next year's homes = this year's homes + completions - demolitions \pm conversions, with real lead times.

9.3 Two Compounding Problems

- Correlation standing in for causation. A statistical or ML model can learn that when incentives
 rose in the past, completions later rose too. But that is a pattern, not a warranted cause; without
 explicit constraints, the model can imply supply that cannot be physically built. Machine learning
 is not "AI making decisions"; it is statistics done quickly and thoroughly. On its own it does not
 respect conservation, lags, or identification unless told to and tested.
- 2. Out-of-distribution drift. Ask an LLM to estimate local impacts or explain mechanisms and, outside well trodden "canonical" narratives, it degrades and diverges confidently giving different numbers and rationales when you vary labour tightness, approvals timing, or servicing constraints. Fluency ≠ fitness-for-purpose.

9.4 WHAT THE TOP-DOWN VIEW IGNORES: THE HOUSING PIPELINE

Delivering a home requires multiple gates that cannot be wished away by averages or historic correlations:

- Labour skills and location. Carpenters, electricians, plumbers, inspectors, and site supervisors are unevenly distributed. Crews cannot teleport; travel radius and union/local rules matter. Available person hours by trade and by zone bound monthly starts and completions.
- Servicing land with linear infrastructure. Water, wastewater, power, gas, and road access must be
 in place before foundations. Trunk capacity and feeder timing create serviced lot bottlenecks;
 upsizing mains or extending a collector road can take seasons to years.



- Approvals and inspections. Rezonings, site plan approvals, building permits, and staged inspections create lead time ladders that cannot be parallelised beyond a point. Their distributions (not just averages) determine cadence.
- Economic events. Interest rate moves, material price shocks (lumber, concrete), lending standards, and developer balance sheets shift feasibility mid stream; some starts stall, some projects cancel.

9.5 BOTTOM-UP ALTERNATIVE: AGENT-BASED MODELLING (ABM) WITH GUARDRAILS

A bottom-up ABM treats households, builders, crews, and lots as agents. Homes move through permit → start → inspection stages → completion; every move creates a vacancy chain elsewhere; demolitions and conversions reduce stock. Labour capacity is tracked by trade × zone, and serviced lot availability is capped by linear infrastructure schedules. This makes stock/flow accounting auditable by construction. Al (including LLMs and ML) is then used as magnitude under tests - to extract tables, draft code scaffolds, format outputs, enumerate scenarios - not to set ends or to assert causal effects. (Farmer & Foley, 2009; Tesfatsion & Judd, 2006; Sterman, 2000)

9.6 DIRECTION FIRST: A TWO-PAGE VECTOR BRIEF (EXPANDED EXCERPT)

- Vision. Estimate net new homes over three years for a mayor's briefing; deliver a four page note with one figure on bottlenecks.
- Ethics. No causal language without identified designs or calibrated simulation; report who benefits/who waits; disclose uncertainty.
- Constraints (Source-lock). Use only official permits/starts/completions/demolitions, utility servicing schedules, inspection logs, trades registry counts by zone, and published material price indices. Any new source requires logged re authorisation.
- Trade offs. Prioritise pipeline realism over speculative price paths; run two capacity cases (baseline labour vs. tight labour) and two servicing cases (on time vs. slippage).
- Objectives (Acceptance tests).
 - 1. Conservation test: $Stock(t+1) = Stock(t) + Completions Demolitions \pm Conversions (tolerance <math>\pm 0.1\%$).
 - 2. Pipeline continuity: Starts in zone $z \le$ serviced lots(z) and \le crew hours(z)/hours per start, by trade; completions bounded by lagged starts and inspection capacity.
 - 3. Approvals realism: Lead time distributions respected (no "instant permits"); inspection calendars honoured.
 - 4. Economic events: Shock scenarios (±100 bps rates; ±15% material prices) change feasibility and timing, not just demand curves.
 - 5. Provenance: Each figure regenerates from code; each claim maps to a source or cell; causality labels applied (Causal/Associative/Speculative).



Rigour & roles. Author ≠ Checker ≠ Approver; decision log records any change to V.E.C.T.O.R; LLM
outputs must cite from the source pack; ML fits may inform calibration but cannot violate
conservation or gating tests.

9.7 MAGNITUDE NEXT: WHAT AI/ML CAN DO SAFELY

- LLM: Parse council minutes and utility PDFs into servicing milestone tables; normalise permit statuses; generate figure captions; draft uncertainty language to a rubric.
- ML/statistics: Fit approval time and inspection time distributions; estimate labour productivity curves by trade; forecast material price bands.
- Human conductor: Choose identification or calibration strategy; set capacity caps and serviced lot gating; adjudicate trade offs; sign the rationale.
- Checker: Verify provenance, data rights, and that acceptance tests (conservation, continuity, approvals realism) pass.

9.8 OUTCOME (WHY THE NUMBERS CHANGE)

Under the rails above, the ABM produces ~2,000–3,500 net new homes in three years - not 10,000. Labour is binding in two zones; serviced lot slippage delays start by a season; approvals queues push completions outside the window; a midperiod rate hike shifts a tranche of projects from feasible to paused. Roughly half of the "uplift" in the top-down slide was reshuffling (pull forward purchases, tenure switches), and some "new supply" was netted out by demolitions/conversions. The memo is auditable and reproducible, findings are labelled Associative (no experiment), and the recommendation targets bottlenecks (inspection staffing; trunk main upgrade) rather than celebrating a phantom surge.

Clarifying the tools:

- Machine learning ≠ autonomous AI. ML is statistics at scale; it discovers patterns but does not enforce conservation or causality unless the brief and tests demand it. (Breiman, 2001)
- LLMs ≠ ground truth. They draft and format quickly but degrade and diverge beyond canon; source lock and labels keep them honest. (Koh et al., 2021; Lin, Hilton, & Evans, 2022)
- Conductor + Vector Brief turn into instruments of magnitude under direction: fast where it's safe, cautious where it matters.

9.9 WHY THIS SUPPORTS THE PAPER'S ARGUMENT

It shows how aggregate or Massed projections, unmoored from pipeline physics, manufacture supply on paper; how LLMs amplify the error when prodded beyond familiar patterns; and how the Conductor Model, codified through a short Vector Brief with causality labelling, source lock, acceptance tests, and role separation, produces a usable vector - direction with magnitude - so decisions can be fast and defensible.



10.0 CONCLUSION: LEADERSHIP WITH LEVERAGE

Abundant model capacity has permanently shifted the production frontier of knowledge work. What has not shifted is accountability for ends, trade-offs, and standards. The central claim of this paper is plain: Al supplies scalable magnitude; humans supply direction. Treat the relationship as a contest - human versus Al - and you either abdicate judgement to pattern completers or refuse tools that could multiply your team's value. Treat it as a division of cognitive labour and the pair forms a usable vector: direction with magnitude.

The practical move, then, is to stop asking who wins and start asking who does what. When we cast the machine as a partner rather than a rival, two traps lose their pull: the urge to let systems choose our goals (automation bias) and the habit of guarding every task as artisan work (manual pride). Clarity about roles turns speed into an asset instead of a liability.

This is where the conductor comes in. Someone accountable must set the rails: purpose, values, constraints, trade offs, measurable objectives, and rigour. That's not extra paperwork; it is the operating envelope that lets magnitude run without eroding judgement. A simple, durable rule follows delegate when inputs, sources, and acceptance tests are explicit, and failure is cheap to catch; retain when a task defines what counts as good or reconciles values and risks; work jointly when the machine can widen options and a human must still synthesise and decide.

Causality needs special care. Modern systems are probabilistic patterns; if we do not tell them where causal language is allowed - and on what evidence - they will produce fluent continuation, not warranted cause and effect. The remedy is straightforward: label claims Causal / Associative / Speculative; lock citations to a permitted source pack; make outputs reproducible by construction; and keep a decision log of what changed and why. Add role separation - Author \neq Checker \neq Approver - and you preserve human command of ends while letting machines accelerate the means.

Used inside these rails, AI shines at what most slows teams down: recall at scale, pattern completion across formats, multi-style drafting, formatting and assembly, code scaffolding, combinatorial exploration, and coordinated summarisation - especially once a schema is set. It stumbles where leadership belongs choosing goals, managing cross-domain trade-offs, making value judgements, exercising epistemic caution, handling theory and identification, surviving out-of-distribution corners, maintaining long-horizon coherence, and ensuring compliance. In plain terms: machines excel at magnitude under tests; humans own direction under standards.

Ignore the conductor and familiar risks scale up fast: accountability drifts, provenance slips, metrics get gamed, causal claims get overstated, scope and norms slide, and throughput becomes irreproducible. The fixes are structural, not heroic - keep roles separate, require source-lock and provenance tables, set acceptance bars before work begins, and make uncertainty statements a condition of release. Speed without these rails multiplies error.

What keeps the system honest is conceptual knowledge. Think of it as the team's loss function: mechanisms, measures, and limits turned into acceptance tests that penalise wrong answers by design.



Without this discipline, models amplify category errors, proxy traps, path dependence, Goodhart effects, and normative slippage. With it, the same tools become disciplined instruments whose speed compounds value rather than compounding error.

To make this usable every day, we offered a compact authorisation: V.E.C.T.O.R. - Vision, Ethics, Constraints, Trade-offs, Objectives, Rigour. You set these before runs and test against them after. The guardrails are baked in causality labels, source-lock, acceptance tests, uncertainty and reproducibility, and role separation. The detailed checklists and templates live in the companion paper; here, the point is the habit: fix direction so magnitude cannot erode accountability.

And the story is not abstract. We walked through housing supply to show how a tidy, top-down projection - sometimes dressed in machine learning correlations - can put "10,000 homes" on a slide by smoothing over stock and flow conservation, trade specific labour capacity, service and timing, approvals ladders, and economic shocks. A bottom up, agent-based approach - run under a short Vector Brief - returned ~2,000–3,500 real, auditable units and pointed straight at the bottlenecks to fix. The lesson generalises: ML is statistics at scale - not autonomous judgement - and LLMs degrade and diverge beyond canon. Only the conductor's guardrails keep both useful.

So, the next steps are plain. Institutionalise the conductor role. Use V.E.C.T.O.R. as the standard brief. Require source lock, provenance, causality labelling, uncertainty, reproducibility, and role separation for every AI assisted artefact. Measure what improves decisions - traceability, uncertainty coverage, decision latency - and train teams to write acceptance tests before analysis and to recognise when the model has wandered out of lane. That is leadership, not lubrication: the human sets direction and standards; the machine multiplies labour inside those rails. (W3C, 2013; Office for Statistics Regulation, 2021)

The conductor is not a bottleneck; the conductor creates usable vectors - pairing direction with magnitude - so organisations can move not merely faster, but on purpose.



A. BIBLIOGRAPHY

Bender, E. M., Gebru, T., McMillan-Major, A., & Shmitchell, S. (2021). On the dangers of stochastic parrots: Can language models be too big? In *Proceedings of the 2021 ACM Conference on Fairness, Accountability, and Transparency (FAccT '21)* (pp. 610–623). https://doi.org/10.1145/3442188.3445922

Breiman, L. (2001). Statistical modeling: The two cultures (with comments and a rejoinder by the author). *Statistical Science*, *16*(3), 199–231. https://doi.org/10.1214/ss/1009213726

Dell'Acqua, F., McFowland III, E., Mollick, E., Lifshitz-Assaf, H., Kellogg, K. C., Rajendran, S., Lakhani, K. R. (2023). *Navigating the jagged technological frontier: Field experimental evidence of the effects of AI on knowledge worker productivity and quality* (HBS Working Paper No. 24-013). European Union. (2024). Regulation (EU) 2024/1689 of the European Parliament and of the Council of 13 June 2024 on artificial intelligence (AI Act). *Official Journal of the European Union*. https://eurlex.europa.eu/eli/reg/2024/1689/oj

Farmer, J. D., & Foley, D. (2009). The economy needs agent-based modelling. *Nature, 460*, 685–686. https://doi.org/10.1038/460685a

Goodhart, C. A. E. (1975). Problems of monetary management: The U.K. experience. In *Papers in Monetary Economics* (Vol. 1). Reserve Bank of Australia.

Hendrycks, D., & Gimpel, K. (2017). A baseline for detecting misclassified and out-of-distribution examples in neural networks. In *Proceedings of the International Conference on Learning Representations (ICLR)*. https://arxiv.org/abs/1610.02136

Hendrycks, D., Basart, S., Mazeika, M., Zou, A., Kwon, J., Mostajabi, M., Steinhardt, J., & Song, D. (2019). Scaling out-of-distribution detection for real-world settings. In *Advances in Neural Information Processing Systems (NeurIPS)*. https://arxiv.org/abs/1911.11132

Hernán, M. A., & Robins, J. M. (2020). *Causal inference: What if.* Boca Raton, FL: Chapman & Hall/CRC. https://doi.org/10.1201/9781315374932

International Energy Agency. (2024). *Energy and AI*. Paris: IEA. https://www.iea.org/reports/energy-and-ai. (IEA)ISO/IEC. (2023). *ISO/IEC 42001:2023—Information technology—Artificial intelligence—Management system*. International Organization for Standardization.

Koh, P. W., Sagawa, S., Marklund, H., Xie, S. M., Zhang, M., Balsubramani, A., Hu, W., Yasunaga, M., Phillips, R. L., Gao, I., Lee, T., David, E., Stavness, I., Guo, W., Earnshaw, B., Haque, I., Beery, S., Leskovec, J., Kundaje, A., ... Liang, P. (2021). WILDS: A benchmark of in-the-wild distribution shifts. In *Proceedings of the 38th International Conference on Machine Learning (ICML)*. https://proceedings.mlr.press/v139/koh21a.html

Lewis, P., Perez, E., Piktus, A., et al. (2020). Retrieval-augmented generation for knowledge-intensive NLP tasks. *NeurIPS 2020*. https://arxiv.org/abs/2005.11401. (NeurIPS Proceedings)



Lin, S., Hilton, J., & Evans, O. (2022). TruthfulQA: Measuring how models mimic human falsehoods. In *Proceedings of the 60th Annual Meeting of the Association for Computational Linguistics (ACL)* (pp. 3214–3252). https://doi.org/10.18653/v1/2022.acl-long.229

NIST. (2023). *Artificial Intelligence Risk Management Framework (AI RMF 1.0)* (NIST AI 100-1). National Institute of Standards and Technology. https://doi.org/10.6028/NIST.AI.100-1

Office for Statistics Regulation. (2021). *Reproducible Analytical Pipelines: Overcoming barriers to adoption*. UK Statistics Authority.

OECD. (2024). *OECD framework for the classification of AI systems*. Organisation for Economic Cooperation and Development.

OWASP Foundation. (2025). OWASP Top 10 for Large Language Model Applications (v2025). https://owasp.org/www-project-top-10-for-large-language-model-applications/ (OWASP Foundation)Pearl, J. (2009). Causality: Models, reasoning, and inference (2nd ed.). Cambridge University Press.

Parasuraman, R., & Riley, V. (1997). Humans and automation: Use, misuse, disuse, abuse. *Human Factors,* 39(2), 230–253. https://doi.org/10.1518/001872097778543886

Sterman, J. D. (2000). *Business dynamics: Systems thinking and modeling for a complex world*. Boston, MA: Irwin/McGraw-Hill.

Strathern, M. (1997). 'Improving ratings': Audit in the British University system. *European Review, 5*(3), 305–321.

Tesfatsion, L., & Judd, K. L. (Eds.). (2006). *Handbook of computational economics, Vol. 2: Agent-based computational economics*. North-Holland.

Vector Institute. (2025). *Principles in action: A playbook for responsible AI product development*. Toronto, Canada: Vector Institute.

