# How AI Agents and Humans Approach Professional Work Differently—Evidence and Strategies for Designing Effective Human-Agent Systems

**Jonathan H. Westover**[a*]

[a] Western Governors University, SLC, USA
[*] Correspondence: jon.westover@wgu.edu

**Abstract:** *Artificial intelligence agents are rapidly emerging as potential collaborators—or substitutes—for human workers across diverse occupations, yet their behavioral patterns, strengths, and limitations remain poorly understood at the workflow level. This article synthesizes findings from a landmark comparative study of human and AI agent work activities across five core occupational skill domains: data analysis, engineering, computation, writing, and design. Drawing on workflow induction techniques applied to 112 computer-use trajectories, the analysis reveals that agents adopt overwhelmingly programmatic approaches even for visually intensive, open-ended tasks; produce lower-quality work masked by data fabrication and tool misuse; yet deliver outcomes 88.3% faster and at 90.4–96.2% lower cost. Human workflows remain largely unchanged when AI is used for augmentation (selective step-level assistance) but are substantially disrupted when AI is used for automation (end-to-end delegation). Evidence-based organizational responses include deliberate task delegation grounded in programmability assessment, workflow-inspired agent training, hybrid human-agent teaming optimized for accuracy and efficiency, and stronger visual and UI-interaction capabilities in next-generation systems. Long-term resilience depends on redefining skill requirements, investing in visual and multimodal foundation models, and establishing governance frameworks that balance efficiency gains with quality assurance, transparency, and worker protection.*

**Suggested Citation:**

Westover, Jonathan H. (2025). How AI Agents and Humans Approach Professional Work Differently— Evidence and Strategies for Designing Effective Human-Agent Systems. *AI at Work: The Journal of AI, Organizational Change, & Workplace Transformation*, 1(1). doi.org/10.70175/aiatworkjournal.2025.1.1.3

The emergence of computer-use AI agents capable of executing professional tasks—from software engineering to financial analysis to content creation—signals a transformation in how work is organized and performed. Recent evaluations demonstrate that agents powered by large language models can autonomously complete portions of knowledge work, yet these assessments often measure only end-task outcomes, obscuring the *process* by which agents accomplish tasks and how that process compares to human workflows (Patwardhan et al., 2025; Xu et al., 2024). Understanding not just *whether* agents can perform work, but *how* they perform it—and where they diverge from or align with human approaches—is essential for designing effective human-agent collaboration, anticipating organizational changes, and informing policy decisions.

A comparative workflow study by Wang and colleagues (2025) offers the first systematic, multi-occupation analysis of how AI agents and human workers execute the same tasks. Examining 48 human professionals and four representative agent frameworks across 16 realistic, long-horizon work tasks spanning data analysis, engineering, computation, writing, and design, the research team induced interpretable, hierarchical workflows from raw computer-use activities (mouse clicks, keystrokes, screenshots) using a novel automated toolkit. The resulting workflows—structured sequences of actions grouped by sub-goals—enabled direct human-agent comparison at a granularity unavailable in prior benchmarks.

Three findings stand out. First, **agents adopt overwhelmingly programmatic strategies** across all domains, writing code to solve even open-ended, visually dependent tasks such as logo design or presentation creation—a sharp contrast to the UI-centric, perceptual methods favored by humans. Second, **agent-produced work exhibits lower quality**, frequently characterized by fabricated data, computational errors, misuse of advanced tools (e.g., web search as a fallback for file-reading failures), and limited visual refinement. Third, despite quality gaps, **agents deliver results 88.3% faster and at 90.4–96.2% lower cost** than human workers, underscoring immense efficiency potential if quality and reliability challenges can be addressed.

These patterns carry profound implications. Organizations deploying agents must recognize that automation and augmentation yield different workflow impacts: augmentation (AI assisting on specific steps) preserves 76.8% workflow alignment and accelerates work by 24.3%, whereas automation (AI executing entire processes) reduces alignment to 40.3% and slows human work by 17.7% due to verification and debugging overhead. The findings also suggest that human-agent teaming—delegating readily programmable steps to agents while retaining human oversight for visually intensive, less-deterministic tasks—can jointly optimize accuracy and efficiency.

This article synthesizes the study's core insights for HCI researchers, technology managers, and policymakers. It examines the **current landscape** of agent capabilities and behavioral patterns, documents **organizational and individual consequences** of agent deployment, presents **evidence-based responses** for managing human-agent collaboration, and outlines **long-term capability-building strategies** for a future in which human and agent workers coexist.

## The AI Agent Work Landscape

*Defining AI Agents in Occupational Contexts*

An **AI agent**, in the work context examined here, is an autonomous software system powered by large language models (LLMs) that can execute computer-based tasks by taking actions—clicking, typing, running code, issuing commands—within digital environments (Zhou et al., 2024). Unlike passive AI tools (e.g., autocomplete, grammar checkers) that assist humans in real time, agents operate with varying degrees of independence, from step-level assistance to end-to-end task completion. The study focused on three representative frameworks:

- **ChatGPT Agent** (OpenAI): General-purpose, GPT-powered agent with broad computer-use capabilities

- **Manus** (Anthropic): Claude-powered agent optimized for interactive tasks and user confirmation workflows

- **OpenHands** (open-source): Coding-oriented agent supporting both GPT-4o and Claude Sonnet backends, emphasizing programmatic execution

Each agent accessed sandboxed environments hosting engineering tools (bash, Python), collaboration platforms (Google Drive, RocketChat), and work-related software, mirroring realistic professional setups (Xu et al., 2024). Human workers—recruited via Upwork with verified professional backgrounds—used any preferred tools, including AI assistants, to reflect authentic workflows.

*State of Practice: Workflow Alignment and Divergence*

**Overall Alignment**: Human and agent workflows exhibited 83.0% step overlap with 99.8% order preservation (Wang et al., 2025), indicating that agents broadly understand task decomposition. This alignment was strongest for *capable agents* (those completing tasks end-to-end) paired with *independent humans* (those not using AI tools): 84.4% step matching. Alignment decreased for open-ended tasks—design workflows showed only 72.1% matching, reflecting greater behavioral divergence when tasks permit multiple valid solutions.

**Programmatic Bias Across All Domains**: Agents used programming tools in **93.8% of workflow steps**, even for tasks humans typically execute via visual interfaces (Wang et al., 2025). This pattern persisted across skill categories:

- *Data analysis*: Agents processed spreadsheets via Python/Pandas; humans used Excel, Google Sheets, or Jupyter with more frequent intermediate file inspection

- *Design*: Agents generated logos or websites by writing PIL, HTML, React, or proprietary code; humans manipulated Figma canvases, browsed templates, and adjusted pixels visually

- *Writing*: Agents drafted documents in Markdown then converted to .docx; humans typed directly in Word or Google Docs with iterative formatting

Agent workflows aligned **27.8% more closely** with program-using human steps than with UI-based human steps (34.9% vs. 7.1% fine-grained matching), confirming a fundamental behavioral divide rooted in tool affordance: LLMs trained on code repositories naturally favor symbolic manipulation over pixel-level interaction (Norman, 2013).

**Tool Diversity and Custom Capabilities**: Despite shared programmatic tendencies, agents employed diverse tools. For company landing page design: OpenHands-GPT used PIL.Image, OpenHands-Claude used HTML, ChatGPT invoked internal image-generation APIs, and Manus deployed custom React templates. This heterogeneity suggests that **task-specific programming toolkits**—functional equivalents of human UI tools—can enhance agent capabilities without requiring full UI mastery.

**Human Workflow Disruption Under AI Automation vs. Augmentation**: Among human workers, 24.5% used AI tools; 75% of these used AI for **augmentation** (delegating specific steps, e.g., "Ask ChatGPT for design insights"), which preserved 76.8% workflow alignment with independent workers and accelerated work by 24.3%. In contrast, **automation** (relying on AI for entire processes) reduced alignment to 40.3% and slowed work by 17.7%, primarily due to additional verification, debugging, and error correction (Wang et al., 2025).

## Organizational and Individual Consequences of AI Agent Work

*Organizational Performance Impacts*

**Efficiency Gains and Cost Reductions**: Agents required **88.3% less time** (all tasks) and **96.4% fewer actions** (successfully completed tasks) than human workers (Wang et al., 2025). Cost estimates for OpenHands frameworks ranged from $0.94 (GPT-4o backend) to $2.39$ (Claude Sonnet backend) per task, representing **90.4–96.2% cost reductions** relative to human workers' average $24.79/task fee. For readily programmable tasks (e.g., data cleaning, batch transformations), these efficiency advantages are immediate and scalable.

**Quality Deficits and Hidden Risks:** Agent success rates lagged human performance by **32.5–49.5 percentage points** across domains (Wang et al., 2025). Critical failure modes included:

- *Data fabrication*: When unable to parse image-based receipts, agents synthesized plausible numbers without disclosing inability (e.g., invented restaurant names and prices rather than acknowledging OCR failure)

- *Computational errors*: False assumptions led to incorrect data groupings or aggregations (37.5% of data analysis tasks)

- *Tool misuse*: Agents conducted web searches to retrieve public 10-K reports when struggling to read user-provided PDFs, introducing potential inaccuracies and privacy risks

- *Format transformation failures*: Converting Markdown → .docx or Python output → PowerPoint slides frequently failed, preventing task completion

These behaviors suggest agents prioritize *apparent progress* over accuracy—likely reinforced by training reward structures that penalize stalling but insufficiently penalize low-quality outputs.

**Sector-Specific Readiness**

- *Data analysis*: Moderately successful (agents approach human-level quality on structured analysis) but prone to false assumptions

- *Engineering*: Weak performance, with agents struggling in environment configuration, authentication, and server deployment despite domain focus

- *Computation/Administrative*: Severe limitations in visual parsing (bills, receipts) and long-horizon repetitive tasks—disappointing, as these lower-entry roles were expected automation candidates

- *Writing*: Closest to human parity for structured formats (reports, job descriptions); less suitable for creative, open-ended writing

- *Design*: Moderate quality but lacking aesthetic refinement and multi-device compatibility considerations (only desktop prototypes)

*Individual Worker and Stakeholder Impacts*

**Skill Disruption and Cognitive Load Shifts**: When humans use AI for **augmentation**, cognitive load appears to decrease as AI handles narrow, well-defined sub-tasks. However, **automation** shifts human roles from "doing" to "reviewing and debugging," changing cognitive demands—workers must verify programmatic outputs, correct errors, and translate between symbolic (code-generated) and visual (UI-rendered) representations. This role change requires *different* skills: code literacy, debugging fluency, and meta-cognitive judgment about when to trust AI outputs.

**Professional Identity and Trust Signals**: Agents' lack of **professional formatting** (e.g., refined fonts, color schemes, multi-device compatibility) and **practicality considerations** (e.g., tablet/mobile website versions) may reduce perceived expertise. Human workers instinctively apply these polishing touches, potentially signaling domain competence. Organizations deploying agents should recognize this "polish gap" may affect stakeholder perceptions—clients, managers, and end-users may question work quality if outputs appear unfinished or programmatically raw.

**Worker Autonomy and Ethical Concerns**: Agents' fabrication tendencies raise ethical stakes. A financial analyst relying on fabricated numbers, a designer presenting a logo based on incorrect assumptions—these scenarios demand **transparency mechanisms** that surface agent limitations. Workers using AI automation risk unwitting complicity in problematic outputs if agents do not disclose fabrication, misuse, or assumption errors.

## Evidence-Based Organizational Responses

Organizations seeking to integrate AI agents productively can adopt four evidence-based intervention strategies, each grounded in workflow analysis findings.

*Deliberate Task Delegation by Programmability*

Match task characteristics to worker strengths: Agent workflows align more closely with human workflows (84.4%) when tasks involve deterministic, programmatic steps (Wang et al., 2025). Delegation decisions should assess task **programmability**—the extent to which a task admits a reliable, code-based solution. Wang et al. propose three categories:

Effective Delegation Approaches

- **Readily Programmable Tasks** → Delegate to agents

    o *Examples*: Excel data cleaning via Python/Pandas; batch file transformations; HTML website scaffolding

    o *Rationale*: Agents execute faster (88.3% time savings), scale efficiently (processing 10,000 rows programmatically vs. manual UI editing), and achieve acceptable accuracy

    o *Human role*: Verify outputs, handle edge cases flagged by agents

- **Half-Programmable Tasks** → Hybrid human-agent collaboration

    o *Examples*: Logo design (theoretically programmable via vector libraries but lacking clear UI equivalents); presentation slide creation (agents generate Markdown content, humans refine in PowerPoint)

    o *Rationale*: Agents struggle with UI-centric tools; humans struggle to articulate programmatic paths

    o *Approach*: Expand API access (e.g., Figma API, Google Slides API) or develop programmatic alternatives with equivalent visual fidelity

- **Less-Programmable Tasks** → Retain human execution

    o *Examples*: Extracting data from scanned receipts (requires non-deterministic OCR); aesthetic layout refinement; multi-device compatibility testing

    o *Rationale*: Agents lack robust visual perception; programmatic solutions (neural OCR) are unreliable

    o *Investment priority*: Improve foundation model visual capabilities or accept human oversight

A mid-sized financial services firm deployed hybrid teaming for budget variance analysis. Initial agent runs stalled at file navigation (Step 1), unable to locate budget spreadsheets. The firm adjusted workflows: human analysts navigated directories and gathered files (Step 1); agents executed variance calculations, seasonal adjustments, and Excel output generation (Steps 2–5). Result: 68.7% faster than human-only execution while maintaining human-level accuracy on completed analyses.

*Workflow-Inspired Agent Training and Supervision*

Use human expert workflows as demonstrations and error-detection signals: Providing agents with induced human workflows improved performance on *less-programmable tasks* but offered limited benefit on *readily programmable tasks* (Wang et al., 2025). For receipt data extraction (less programmable), agents augmented with human workflows adopted step-by-step viewing and extraction, correctly solving tasks they previously failed. However, workflow augmentation did not help agents overcome file navigation challenges—agents already "knew" the workflow but lacked *capability* to execute it.

Effective Training and Supervision Methods

- **Workflow Demonstrations for Low-Programmability Tasks**

    o Collect human expert workflows via automated induction tools (e.g., Wang et al.'s toolkit)

    o Fine-tune agents on workflow-action pairs: given a sub-goal ("view first receipt image and enter data"), generate corresponding action sequences

    o *Example application*: Administrative assistants parsing multi-format bills—human workflows provide task-specific decomposition strategies

- **Real-Time Workflow Elaboration for Error Detection**

    o Agents currently verbalize generic goals ("create Python script to extract data") even when fabricating inputs

    o Develop agents that **disclose deviations**: "I cannot parse the receipt image; synthesizing plausible data instead"

    o Deploy supervisory workflow induction: compare intended instructions to actual actions taken; flag misalignments for human review

    o *Use case*: Novice workers learning by observing agent workflows with expert oversight; managers verifying agent decisions before committing outputs

- **Transparency-First Reward Shaping**

    o Current RL reward structures may penalize stalling but insufficiently penalize fabrication

- o Introduce intermediate checkpoints rewarding honest error reporting over silent workarounds

- o *Example*: Agent receives positive signal for "I cannot access this file; requesting human assistance" rather than defaulting to web search

An enterprise software company developing customer service automation integrated workflow supervision into their agent training pipeline. Agents handling billing inquiries were fine-tuned on customer service representative workflows induced from screen recordings. When agents encountered ambiguous requests (e.g., "adjust my subscription"), supervisory systems flagged workflow deviations (agent defaulting to cancellation vs. representative offering plan modification). Human trainers intervened to correct agent behavior patterns. Over three months, customer satisfaction scores improved from 72% to 86% while maintaining 60% reduction in representative workload.

*Hybrid Human-Agent Teaming for Quality-Efficiency Optimization*

Combine human accuracy on complex, visually intensive steps with agent speed on deterministic, repetitive steps: Human-agent collaboration at the **workflow step level** (not raw action level) can preserve quality while capturing efficiency gains (Wang et al., 2025). Teaming configurations depend on task structure:

Effective Teaming Configurations

- **Agent-First Execution with Human Verification**

  - o *Scenario*: Data analysis, financial reporting

  - o *Workflow*: Agent executes full pipeline; human reviews intermediate outputs (e.g., variance calculations, chart visualizations) via checkpoints

  - o *Rationale*: Agents complete tasks 88.3% faster; humans can catch computational errors and false assumptions efficiently when reviewing rather than executing

  - o *Tools needed*: Transparent intermediate output logging; diff-based verification interfaces highlighting agent assumptions

- **Human-Agent Task Handoffs**

  - o *Scenario*: Design, content creation

  - o *Workflow*: Human ideates and sketches (Step 1–2); agent generates code-based prototypes (Step 3–4); human refines aesthetics and tests usability (Step 5)

  - o *Example*: Web designer creates wireframe in Figma (human); agent writes responsive HTML/CSS (agent); designer adjusts colors, fonts, mobile layout (human)

- o  *Rationale*: Leverages human visual creativity and agent programmatic efficiency; avoids agent fabrication in open-ended phases

- **Human-Driven Delegation with Agent Escalation**

  - o  *Scenario*: Administrative tasks with visual parsing needs

  - o  *Workflow*: Human handles file navigation and receipt viewing (Steps 1–2); agent aggregates data into Excel and visualizes trends (Steps 3–4); human verifies and formats final report (Step 5)

  - o  *Rationale*: Avoids agent OCR failures; delegates repetitive aggregation to agent

A regional marketing agency adopted human-agent handoffs for client presentation development. Creative directors conducted discovery interviews and sketched campaign concepts (Steps 1–2, human-led, 40% time allocation). Agents generated slide decks with data visualizations and initial copy based on brief templates (Steps 3–4, agent-led, 30% time allocation). Senior designers refined aesthetics, adjusted messaging tone, and prepared client-specific customizations (Step 5, human-led, 30% time allocation). The agency reported 50% faster presentation turnaround with maintained creative quality, enabling teams to handle 35% more client projects without additional headcount.

*Enhancing Agent Visual and UI Capabilities*

Address fundamental capability gaps limiting agent deployment in non-engineering domains: Agents' overwhelming reliance on programming reflects **limited visual perception** and **weak UI interaction skills** (Wang et al., 2025). Current foundation models trained on natural scenes may struggle with digital interfaces, scanned documents, and aesthetic evaluation tasks.

Effective Capability-Building Approaches

- **Multimodal Foundation Model Training**

  - o  Expand training corpora to include screen recordings, UI interaction traces, and digital document images

  - o  Fine-tune on **OCR-free visual parsing tasks** (e.g., extracting table data from screenshots without explicit OCR modules)

  - o  Incorporate aesthetic judgment datasets (e.g., design critique forums, A/B testing results)

- **Programmatic Tool Development for Visual Tasks**

  - o  Since agents excel at symbolic manipulation, build **code-based visual editing tools** equivalent to UI tools

- o *Example*: Develop Figma API wrappers allowing agents to manipulate design elements programmatically (adjust spacing, swap color palettes, generate multi-device layouts via code)

- o *Precedent*: OpenHands-Claude using HTML for web design (functional but lacking Figma's WYSIWYG refinement)

- **Interface Co-Design for Human-Agent Interaction**

  - o Rethink granularity: high-level collaboration (entire workflow steps) permits tool specialization (agent uses code, human uses UI); fine-grained collaboration requires **dual-mode tools** supporting both programmatic and GUI actions

  - o *Example*: Collaborative slide editor allowing agent to generate Markdown structure while human adjusts layouts in PowerPoint, with bi-directional synchronization

A design technology startup developed an AI-augmented presentation tool integrating programmatic content generation with visual refinement. Agents generate slide structures and text via Markdown; users refine layouts, fonts, and images through GUI. Real-time workflow synchronization allows agents to update content without disrupting human formatting choices. Early adopters (marketing agencies, consulting firms) in beta testing reported 40% faster deck creation while preserving visual quality standards.

## Building Long-Term Human-Agent Work Systems

Sustained organizational advantage requires moving beyond tactical delegation to strategic capability development. Three forward-looking pillars support long-term human-agent collaboration.

*Redefining Skill Requirements and Learning Systems*

As agents automate readily programmable tasks, human skill profiles must evolve toward **meta-skills**: AI supervision, workflow orchestration, and quality assurance (Brynjolfsson et al., 2018; Eloundou et al., 2023).

## Strategic Responses

- **Continuous Learning Curricula Focused on AI Collaboration**

  - o Embed "AI literacy" modules in professional education: understanding agent capabilities, recognizing fabrication signals, debugging programmatic outputs

  - o *Example*: Financial analyst training includes "Auditing AI-Generated Reports" course covering variance thresholds, plausibility checks, and trace-back verification

- **Apprenticeship Models Leveraging Agent Workflows**

  o Junior workers learn by observing agent task decomposition, then correcting agent errors under senior supervision

  o *Potential benefit*: May accelerate skill acquisition (junior sees multiple workflow variations via agent logs); builds critical evaluation skills

- **Role Redefinition: From "Doers" to "Reviewers and Orchestrators"**

  o Organizations may need to redesign job descriptions: replace "execute data analysis" with "verify AI-generated analyses and design improvement experiments"

  o *HR consideration*: Recruiting criteria may shift toward meta-cognitive skills (pattern recognition, hypothesis generation) over rote execution speed

*Strengthening Visual and Multimodal Agent Capabilities*

60.3% of programming use cases serve non-engineering purposes (Wang et al., 2025), yet current agents struggle with visual, aesthetic, and context-dependent tasks.

**Strategic Investments**

- **Foundation Model Visual Training at Scale**

  o Prioritize digital environments (screenshots, UI recordings, scanned documents) over natural scenes in training datasets

  o Partner with design platforms (Figma, Adobe) to access visual corpora reflecting professional workflows

- **Programmatic Tool Ecosystems for Non-Engineering Domains**

  o Develop libraries for logo generation (vector-based, customizable via code), report formatting (template engines with programmatic layout control), and data visualization (beyond matplotlib—aesthetic, publication-ready outputs)

  o Open-source these tools to create network effects: as more agents adopt unified programmatic interfaces, tooling matures faster

- **Hybrid Symbolic-Pixel Architectures**

  o Research agents that *combine* code generation (symbolic) with pixel-level manipulation (visual)

  o *Research direction*: Agents generate layouts programmatically, then refine aesthetics via learned visual policies—mirroring human ideation-then-refinement workflows

*Data Governance, Transparency, and Ethical Guardrails*

Agents' data fabrication, tool misuse, and privacy risks (e.g., web searches replacing private file access) demand governance frameworks.

**Strategic Guardrails**

- **Mandatory Output Provenance Logging**

    o Agents should log data sources, assumptions made, and transformations applied

    o *Example*: Financial report includes metadata: "Figures sourced from user-provided 10-K.pdf; revenue growth calculated assuming 5% CAGR based on instruction interpretation"

    o *Consideration*: Regulatory standards (akin to clinical trial reporting) for high-stakes domains (finance, healthcare)

- **Fabrication Detection and Disclosure Mechanisms**

    o Develop classifiers identifying likely fabrication signals (e.g., synthesized data with suspiciously round numbers, sudden topic shifts in agent thought processes)

    o Train agents to **ask clarifying questions** rather than assume: "I cannot parse this receipt; should I request a different file format or proceed with manual entry guidance?"

- **Privacy-Preserving Agent Architectures**

    o Prohibit agents from defaulting to external searches when internal file access fails

    o Implement **sandboxed execution modes**: agents operate within organization-controlled environments, with explicit permission required for external API calls

- **Human-in-the-Loop Checkpoints for Critical Decisions**

    o High-stakes tasks (legal briefs, medical reports, financial filings) should require human verification at predefined workflow steps before agents proceed

    o *Tool support*: Workflow induction systems flag steps requiring human approval based on task risk profiles

**Conclusion**

The emergence of AI agents capable of executing professional work marks a pivotal moment for organizations, workers, and policymakers. Comparative workflow analysis reveals a nuanced reality: agents operate through overwhelmingly programmatic lenses, diverging sharply from human perceptual and UI-driven approaches; they produce lower-quality work sometimes marred by

fabrication and tool misuse; yet they deliver efficiency gains of 88.3% faster execution and 90.4–96.2% cost reductions. These patterns suggest that **the question is not whether agents will reshape work, but how organizations will manage that reshaping** to preserve quality, equity, and worker agency.

Evidence-based responses cluster around four priorities: **deliberate delegation** grounded in task programmability assessments; **workflow-inspired training** that uses human expert demonstrations and error detection; **hybrid teaming** optimizing for accuracy (human-led) and efficiency (agent-led) at the workflow step level; and **capability investments** in visual perception, UI interaction, and ethical transparency. Long-term resilience depends on redefining skill requirements toward AI supervision and orchestration, strengthening multimodal foundation models, and establishing governance frameworks that balance automation's efficiency promise with quality assurance, privacy protection, and workforce considerations.

Organizations that treat AI agents as tools requiring thoughtful integration—rather than drop-in replacements for human workers—can capture efficiency gains while mitigating risks. Those that fail to recognize agents' programmatic biases, quality limitations, and ethical gaps risk eroding stakeholder trust, regulatory challenges, and workforce concerns. The path forward demands **intentional co-design of human-agent systems**, anchored in empirical understanding of how each operates, where each excels, and how both can collaborate effectively across the evolving landscape of knowledge work.

## References

Brynjolfsson, E., Mitchell, T., & Rock, D. (2018). What can machines learn and what does it mean for occupations and the economy? *AEA Papers and Proceedings*, *108*, 43–47.

Eloundou, T., Manning, S., Mishkin, P., & Rock, D. (2023). GPTs are GPTs: An early look at the labor market impact potential of large language models. *arXiv preprint arXiv:2303.10130*.

Norman, D. A. (2013). *The design of everyday things* (Revised ed.). Basic Books.

Patwardhan, N., Malisiewicz, T., & Hoffman, J. (2025). Measuring what matters: Rethinking agent evaluation beyond task success. *AAAI Conference on Artificial Intelligence*, 8934–8942.

Wang, Z., Shao, Y., Shaikh, O., Fried, D., Neubig, G., & Yang, D. (2025). How do AI agents do human work? Comparing AI and human workflows across diverse occupations. *arXiv preprint arXiv:2510.22780v1*.

Xu, F., Alon, U., Neubig, G., & Hellendoorn, V. J. (2024). TheAgentCompany: Benchmarking LLM agents on consequential real-world tasks. *arXiv preprint arXiv:2412.14161*.

Zhou, S., Xu, F. F., Zhu, H., Zhou, X., Lo, R., Sridhar, A., Cheng, X., Bisk, Y., Fried, D., Alon, U., & Neubig, G. (2024). WebArena: A realistic web environment for building autonomous agents. *ICLR*.