# Chapter 3. User Experience Design for Agentic Systems

As agent systems become an integral part of our digital environments—whether through chatbots, virtual assistants, or fully autonomous workflows—the user experience (UX) they deliver plays a pivotal role in their success. While foundation models and agent architectures enable remarkable technical capabilities, how users interact with these agents ultimately determines their effectiveness, trustworthiness, and adoption. A well-designed agent experience not only empowers users but also builds confidence, minimizes frustration, and ensures clarity in agent capabilities and limitations. The field of agent UX is evolving at an unprecedented pace. New interface paradigms, modality combinations, and user interaction models are emerging almost monthly. This chapter provides foundational design principles that remain relevant even as the specific technologies and capabilities continue to advance rapidly. Designing UX for agent systems introduces unique challenges and opportunities. Agents can interact through a variety of modalities, including text, graphical interfaces, speech, and even video.

Table 3-1. Placeholder

| Modality | Prevalence | Example use cases | Ideal situations |
|---|---|---|---|
| Text | Very common | Customer service chatbots, productivity assistants | When clear, asynchronous, or searchable communication is needed |
| Graphical user interfaces (GUI) | Common | Workflow orchestration dashboards, AI coding assistants like Cursor | When visual structure, context management, or multistep workflows are important |
| Speech/voice | Less common | Siri, smart home assistants (Alexa, Google Home), call center automation | When hands-free interaction or natural conversation is required |
| Video | Rare | Virtual tutors, therapy avatars, interactive learning agents | When visual demonstration, rich expression, or immersive learning is needed |

Another key UX consideration is how the context is managed over time. Some generative AI applications have no memory or learning, so have precisely the information you present them with in exactly that session. This requires users to copy and paste information into the prompt. More modern applications automatically manage this context for you. For example, Cursor uses the integrated development environment to intelligently identify code to include in each model inference. Some applications retain memory over time, enabling agents to remember past interactions, maintain conversation flow, and adapt to user preferences over time. Without these capabilities, even technically advanced agents risk feeling disjointed or unresponsive. Similarly, communicating agent capabilities, limitations, and uncertainty is essential for setting realistic user expectations and preventing misunderstandings. Users must know what an agent can and cannot do, and when they might need to intervene or provide guidance.

Finally, trust and transparency remain foundational to positive user experiences with agent systems. Predictable agent behavior and clear explanations of actions contribute to building relationships where users feel confident relying on agents in high-stakes scenarios.

This chapter explores these core aspects of UX design for agentic systems, offering principles, best practices, and actionable insights to help you design interactions that are intuitive, reliable, and aligned with user needs. Whether you're building a chatbot, an AI-powered personal assistant, or a fully autonomous workflow agent, the principles in this chapter will help you create meaningful and effective experiences that users can trust.

# Interaction Modalities

Agent systems interact with users through a variety of modalities, each offering unique strengths, limitations, and design considerations. Whether through text, graphical interfaces, speech, or video, the choice of modality shapes how users perceive and interact with agents. Text-based interfaces excel in clarity and traceability; graphical interfaces offer visual richness and intuitive controls; voice interactions provide hands-free convenience; and video interfaces enable dynamic, real-time communication.

In the next section, we'll explore these interaction modalities, examining their key strengths, challenges, and best practices for delivering exceptional UX in agent systems.

## Text-Based

Text-based interfaces are one of the most common and versatile ways users interact with agent systems—found in everything from customer service chatbots and command-line tools to productivity assistants integrated into messaging platforms. Their widespread adoption can be attributed to their simplicity, familiarity, and ease of integration into existing workflows. Text interfaces offer a unique advantage: they can support both synchronous conversations (in real time) and asynchronous interactions (where users can return to the conversation at their convenience without losing context). Additionally, text interactions create a clear and traceable record of exchanges, enabling transparency, accountability, and easier troubleshooting when something goes wrong.

In recent years, the text-based modality has undergone a renaissance driven by the integration of advanced AI capabilities within terminal environments. Tools like Warp, Claude Code, and Gemini CLI illustrate this shift vividly. Warp reimagines the traditional developer terminal by integrating natural language command translation, intelligent autocompletion, and context-aware explanations, turning the command line into a collaborative, AI-augmented workspace. To illustrate this trend, Figure 3-1 shows an example of an AI-enabled terminal interface inspired by modern tools like Claude Code and Gemini CLI. This demonstration captures how developers can interact with the terminal using natural language prompts to generate, run, and debug commands seamlessly, without memorizing complex syntax or flags.

Similarly, Claude Code and Gemini CLI extend natural language interactions to code generation, execution, and file manipulation directly within terminal workflows, enabling developers to perform complex tasks by simply describing their goals in plain English. This figure highlights how AI is revitalizing the humble terminal, transforming it from a tool accessible only to those with deep command-line expertise into an approachable, powerful gateway for both novice and expert users to interact with systems through natural language.

This trend reflects a broader rethinking of what text-based interfaces can achieve. The incredible natural language understanding capabilities of modern foundation models are making ordinary text-based interactions more powerful than ever before. Where traditional terminals required precise syntax knowledge and memorization of command flags, AI terminals now act as conversational partners, interpreting user intent, suggesting best practices, and even debugging errors in real time. This shift is democratizing access to powerful systems operations, scripting, and data workflows, making the terminal "new again" as an accessible, intelligent gateway for both novice and expert users.

Figure 3-1. AI-enabled terminal interface. A demonstration of an AI-augmented terminal, where natural language inputs are interpreted into executable commands. Such interfaces transform the traditional command line into an intelligent conversational partner for system operations and development workflows.

However, a key limitation of text-based interfaces is discoverability. Users often do not know what capabilities the agent supports or how to phrase commands effectively. Unlike graphical interfaces—where options, buttons, and menus visually indicate what actions are possible—text-based interfaces require users to guess or recall available functionalities. This lack of affordances can lead to confusion, underutilization of agent capabilities, and user frustration when their requests fall outside the agent's supported scope. For example, a user might ask a support chatbot to modify an order detail that the system does not support, receiving an opaque rejection rather than guidance toward what is possible.

Designing effective text-based agents therefore requires strategies to enhance discoverability. Agents should proactively communicate their supported functions, either through onboarding messages, periodic capability reminders, or dynamic suggestions during conversation. For instance, an agent can respond to a greeting not only with "How can I help you today?" but also with "I can help you cancel orders, check delivery status, or update your account details." This approach ensures users understand the agent's operational boundaries, reducing trial-and-error interactions.

Beyond discoverability, text-based design requires careful attention to clarity, context retention, and error management. Agents should communicate with concise and unambiguous responses, avoiding overly technical jargon or long-winded explanations that may overwhelm the user. Maintaining context across multiturn conversations is equally important; users should not need to repeat themselves or clarify past instructions. Effective agents are also graceful in failure, providing clear error messages and fallback mechanisms, such as escalating to a human operator or offering alternative suggestions when they cannot fulfill a request. Turn-taking management is another subtle but crucial element—agents must guide conversations naturally, balancing when to ask follow-up questions and when to pause for user input.

Ambiguity in natural language remains a significant hurdle, as users may phrase requests in unexpected ways, requiring robust intent recognition to avoid misunderstandings. Additionally, text-based agents are often constrained by response length limits—too short, and they risk being cryptic; too long, and they risk overwhelming or frustrating the user. Emotional nuance is another limitation. Without vocal tone, facial expressions, or visual cues, text-based agents must rely on carefully crafted language to ensure they convey empathy, friendliness, or urgency where appropriate.

Despite these challenges, text-based agents shine in scenarios where precision, traceability, and asynchronous communication are valuable. They excel in customer support, where chatbots provide quick answers to frequently asked questions, or in productivity tools, where command-line interfaces help users execute tasks efficiently. They are equally effective in knowledge retrieval systems, answering specific questions or pulling data from structured databases.

When designed thoughtfully, text-based agents are reliable, adaptable, and deeply useful across a wide range of contexts. For example, text-based agents might be ideal for chat interfaces over messaging apps—like Slack, Teams, and WhatsApp for scalable communications with customers or employees—or text-heavy workloads like customer service, claims processing, or textual research tasks. Their accessibility and ease of deployment make them a cornerstone of agentic UX design—provided their limitations (particularly around discoverability) are mitigated through clear communication of capabilities, robust error handling, and a focus on seamless conversational flow.

# Graphical Interfaces

Graphical interfaces offer users a visual and interactive way to engage with agent systems, combining text, buttons, icons, and other graphical elements to facilitate communication. These interfaces are particularly effective for tasks requiring visual clarity, structured workflows, or multistep processes, where pure text or voice interactions may fall short. Common examples include dashboard-based AI tools, graphical chat interfaces, and agent-powered productivity platforms with clickable elements.

The key strength of graphical interfaces lies in their ability to present information visually and reduce cognitive load. Humans primarily rely on visual input and can process visual information more quickly and easily than text-based information. Well-designed interfaces can display complex data, status updates, or task progress in an intuitive and digestible format. Visual cues, such as progress bars, color coding, and alert icons, guide users effectively without requiring lengthy explanations.

For example, an agent managing a workflow might use a dashboard to show pending tasks, completed steps, and error notifications, enabling users to quickly understand the system's state at a glance. Tools like LangSmith, n8n, Arize, and AutoGen are beginning to illustrate agent workflows visually, making them easier to understand, debug, and reason about; we are likely to see much more of this visual orchestration in the future. To see how these graphical orchestration interfaces are emerging in practice, Figure 3-2 shows an example of a modern agent workflow builder. Tools like this illustrate agent actions, tool calls, conditionals, and outputs as connected visual nodes, enabling developers and operators to easily understand, debug, and optimize complex agentic flows without stepping through raw code alone.
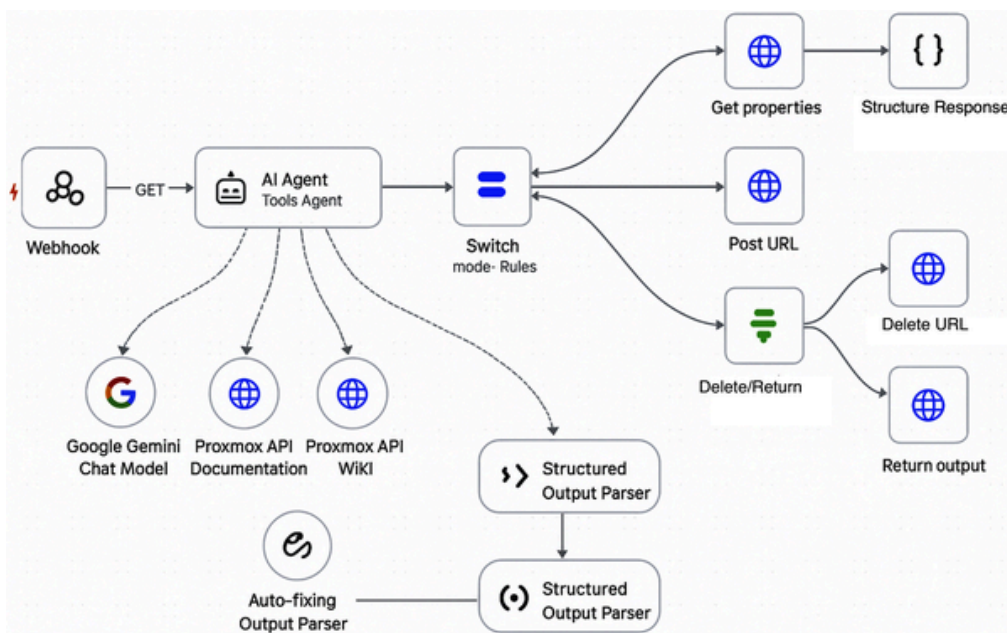
Figure 3-2. Visual orchestration of an agent workflow in n8n.io. This interface displays an AI agent integrated with multiple tools, models, and structured parsing components arranged in a node-based workflow. Such visual designs make it easier to build, manage, and iterate on multistep agent pipelines at scale.

Similarly, Figure 3-3 shows a modern AI-enabled IDE interface, similar to tools like Cursor, Windsurf, Cline, and many more. These environments integrate natural language understanding directly into the coding workflow, enabling developers to ask questions, generate code, refactor functions, and receive explanations or performance optimizations—all within a single, streamlined graphical interface.
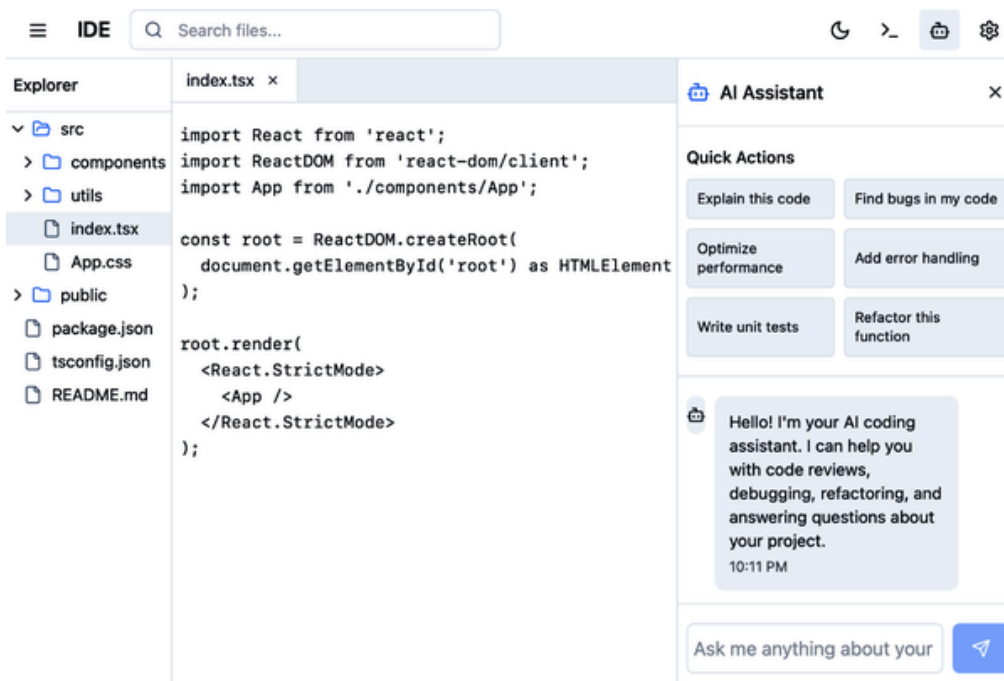


Figure 3-3. AI-enabled IDE interface. An integrated development environment (IDE) enhanced with AI capabilities, combining traditional file explorers and code editors with natural language assistant panels that provide explanations, debugging suggestions, and autogenerated code improvements.

Together, these examples illustrate the rapid evolution of graphical agentic UX. As these interfaces mature, they will redefine what productive, AI-

enabled tools look like—not just for developers, but for every knowledge-intensive profession.

A growing frontier in graphical agent interfaces is the emergence of generative UIs. Instead of relying solely on static dashboards or predesigned layouts, generative UIs dynamically create interface elements, data visualizations, or structured outputs based on user queries. For example, Perplexity AI not only provides textual answers but also generates structured knowledge cards, reference lists, and data tables tailored to the question asked. Similarly, AI coding copilots generate entire forms, config files, or UI components based on user intent.

Generative UIs combine the flexibility of natural language with the clarity and discoverability of graphical layouts, enabling agents to create rich, context-specific interfaces on demand. This expands the usefulness of graphical agents from predefined workflows to open-ended tasks where visual structuring enhances understanding. However, designing generative UIs introduces new challenges: ensuring the generated elements are usable and aesthetically coherent, and that they do not overwhelm users with poorly organized or excessive information. Careful design patterns, layout constraints, and prioritization logic are critical to keep generative UIs effective and user-friendly.

Designing effective graphical agent interfaces also comes with traditional challenges. Screen real estate is limited, requiring prioritization of displayed information to ensure critical details are not buried in clutter. Agents must manage interface responsiveness—users expect real-time updates and smooth transitions between states, especially when agents operate asynchronously. Additionally, graphical elements must adapt gracefully across devices and screen sizes, ensuring consistency whether viewed on a desktop, tablet, or mobile phone.

Another critical consideration is the balance between automation and user control. Graphical interfaces often blend agent autonomy with user-driven actions, such as approving agent-suggested decisions or manually overriding recommendations. For example, an agent suggesting a calendar change might display multiple options through buttons, giving users a clear and efficient way to make a final decision.

Graphical interfaces excel in use cases where data visualization, structured interactions, and clear status updates are essential. Examples include task management dashboards, data analytics tools powered by AI agents, ecommerce product recommendation systems with filters and visual previews, and generative UI systems that dynamically produce structured outputs tailored to user questions. They are particularly effective in hybrid workflows where agents operate in the background but present updates or options visually for user confirmation.

When implemented thoughtfully, graphical and generative interfaces enable clear, efficient, and satisfying interactions with agents. They reduce ambiguity, improve task clarity, and offer users a tangible sense of control. By focusing on clarity, responsiveness, intuitive design patterns, and the emerging potential of generative UI capabilities, graphical interfaces ensure that agent interactions feel smooth, transparent, and aligned with user expectations.

Graphical interfaces excel in use cases where data visualization, structured interactions, and clear status updates are essential. Recent years have seen enormous growth in tools like Lovable, Cursor, Windsurf, and GitHub Copilot, which offer high-quality GUIs that manage context and complex multistep operations with remarkable fluidity. These tools are redefining what productive, agent-enabled interfaces look like for developers. It is time to think just as hard about what the next generation of AI-enabled, agentic UX will be for other professions—lawyers, accountants, insurance professionals, product managers, and knowledge workers. The future of work may not revolve around documents, spreadsheets, and slide decks, but around interactive, agent-driven interfaces purpose-built for decision making, analysis, and creation.

## Speech and Voice Interfaces

Speech and voice interfaces offer users a natural and hands-free way to interact with agent systems, leveraging spoken language as the primary mode of communication. From virtual assistants like Amazon's Alexa and Apple's Siri to customer service voice bots, these interfaces excel in scenarios where manual input is impractical or impossible—such as while driving, cooking, or operating machinery. They also provide an accessible option for users with visual impairments or limited mobility, making agent systems more inclusive.

Historically, latency has been a major barrier for speech and voice interfaces. Processing spoken language in real time—including transcribing speech, interpreting intent, and generating appropriate responses—often led to delays that disrupted conversational flow and made voice interfaces feel clunky or robotic. However, the past two years have seen astonishing advances in this space. New low-latency speech recognition models, combined with more efficient language processing architectures, have dramatically reduced delays. Equally important, the *fluidity* and *capability* of voice AI systems have improved, enabling more natural-sounding interactions that can handle interruptions, mid-sentence corrections, and shifts in conversation topic.

Graceful handling of interruptions is a particularly important aspect of voice interface design. Human conversations are rarely linear monologues; people interrupt themselves to clarify, change direction, or refine a request mid-sentence. Effective voice agents must mirror this conversational flexibility, allowing users to interrupt commands without confusion, revise their inputs seamlessly, and resume where they left off without forcing a complete restart. For example, a user might say, "Book me a table for—oh wait, make that tomorrow instead," and a well-designed agent will adapt fluidly to incorporate the correction without requiring the user to start the command again. This capability not only makes interactions feel more natural but also builds trust and reduces frustration, as users feel the agent is responsive to their real communication patterns rather than demanding rigid, computer-like inputs.

Another major leap has been the integration of tool use into voice agent workflows. Modern voice agents are no longer limited to parsing commands and returning static answers. Instead, they can now pull in external context, update records, and take real-time actions—such as scheduling appointments, changing system configurations, or placing orders—based on dynamic conversational inputs. This ability to combine natural voice interaction with structured backend operations is transforming what voice agents can achieve.

Despite these impressive technological advances, it is important to note that voice interfaces remain a frontier technology. It is true that they have entered mainstream use in smart speakers and simple assistants. However, fully conversational, multiturn, context-aware voice agents with action-taking capabilities are not yet widely deployed across industries. Many enterprises are only beginning to explore voice interfaces for customer service, healthcare, logistics, and field operations.

A key consideration in deploying voice interfaces is understanding the speed at which humans process spoken versus written information. Humans typically speak at 150–180 words per minute, whereas reading speeds average 250–300 words per minute, with skimming speeds exceeding 500 words per minute. This means spoken interfaces are inherently slower for dense or complex information, where text-based interfaces enable faster comprehension and easier reference. However, voice excels in scenarios where hands-free convenience, natural interaction, and immediate contextual responsiveness outweigh these speed constraints.

The following example demonstrates a minimal FastAPI server using the OpenAI Realtime Voice API. It streams microphone audio from a browser to the agent and plays back the assistant's audio responses in real time. Notably, it handles interruptions gracefully: if the user starts speaking mid-response, it immediately truncates the assistant's output to keep the conversation natural. This compact implementation shows the core architecture for building low-latency, interruption-aware voice interfaces with agents:

```python
import os, json, base64, asyncio, websockets
from fastapi import FastAPI, WebSocket
from dotenv import load_dotenv

load_dotenv()
OPENAI_API_KEY = os.getenv("OPENAI_API_KEY")
VOICE          = "alloy"                    # GPT-4o voice
PCM_SR         = 16000                      # sample-rate
PORT           = 5050

app = FastAPI()

@app.websocket("/voice")
async def voice_bridge(ws: WebSocket) -> None:
    """
    1. Browser opens ws://host:5050/voice
    2. Browser streams base64-encoded 16-bit mono PCM c
    3. We forward chunks to OpenAI Realtime (`input_auc
    4. We relay assistant audio deltas back to the brow
    5. We listen for 'speech_started' events and send a
       user interrupts
    """
    await ws.accept()

    openai_ws = await websockets.connect(
```

```python
        "wss://api.openai.com/v1/realtime?" +
            "model=gpt-4o-realtime-preview-2024-10-01"
        extra_headers={
            "Authorization": f"Bearer {OPENAI_API_KEY}'
            "OpenAI-Beta" : "realtime=v1"
        },
        max_size=None, max_queue=None # unbounded for d
    )

    # initialize the realtime session
    await openai_ws.send(json.dumps({
        "type": "session.update",
        "session": {
            "turn_detection": {"type": "server_vad"},
            "input_audio_format": f"pcm_{PCM_SR}",
            "output_audio_format": f"pcm_{PCM_SR}",
            "voice": VOICE,
            "modalities": ["audio"],
            "instructions": "You are a concise AI assis
        }
    }))

    last_assistant_item = None        # track current
    latest_pcm_ts         = 0          # ms timestamp
    pending_marks         = []

    async def from_client() -> None:
        """Relay microphone PCM chunks from browser → C
        nonlocal latest_pcm_ts
        async for msg in ws.iter_text():
            data = json.loads(msg)
            pcm = base64.b64decode(data["audio"])
            latest_pcm_ts += int(len(pcm) / (PCM_SR * 2
            await openai_ws.send(json.dumps({
                "type": "input_audio_buffer.append",
                "audio": base64.b64encode(pcm).decode('
            }))

    async def to_client() -> None:
        """Relay assistant audio + handle interruptions
        nonlocal last_assistant_item, pending_marks
        async for raw in openai_ws:
            msg = json.loads(raw)

            # assistant speaks
            if msg["type"] == "response.audio.delta":
```

```python
                pcm = base64.b64decode(msg["delta"])
                await ws.send_json({"audio":
                    base64.b64encode(pcm).decode("ascii
                last_assistant_item = msg.get("item_id'

            # user started talking → cancel assistant s
            started = "input_audio_buffer.speech_starte
            if msg["type"] == started and last_assistar
                await openai_ws.send(json.dumps({
                    "type": "conversation.item.truncate
                    "item_id": last_assistant_item,
                    "content_index": 0,
                    "audio_end_ms": 0    # stop immediat
                }))
                last_assistant_item = None
                pending_marks.clear()

    try:
        await asyncio.gather(from_client(), to_client()
    finally:
        await openai_ws.close()
        await ws.close()

if __name__ == "__main__":
    import uvicorn
    uvicorn.run("realtime_voice_minimal:app", host="0.0
```

Looking ahead, we are likely to see significant adoption of advanced voice interfaces in the coming years, driven by falling costs, reduced latency, improved speech recognition, and better orchestration with backend tools. In healthcare, voice agents can assist doctors with hands-free note-taking during patient consultations. In customer service, they are replacing rigid interactive voice response (IVR) systems with fluid, humanlike conversations that resolve issues end to end. In industrial applications, workers can control machinery, log observations, or access manuals without stopping their tasks.

Ultimately, voice interfaces are most effective for short, hands-free tasks, quick queries, and action-oriented workflows, rather than for dense information consumption or complex decision making that requires rapid skimming or side-by-side comparison.

When thoughtfully designed, speech and voice interfaces offer unparalleled convenience, accessibility, and flexibility in agent interactions. As these

technologies continue to mature and integrate deeply with backend tools and knowledge systems, they are poised to become indispensable in daily workflows, personal assistants, and enterprise solutions—fundamentally transforming how users interact with AI-powered agents.

## Video-Based Interfaces

Video-based interfaces are an emerging modality for agent interactions, blending visual, auditory, and sometimes textual elements into a single cohesive experience. These interfaces can range from video avatars that simulate face-to-face conversations to agents embedded in real-time video collaboration tools. As video becomes more pervasive in our digital lives—through platforms like Zoom, Microsoft Teams, and virtual event spaces—agents are finding new ways to integrate into these environments. While many of these experiences are still in the uncanny valley, the rapid pace of improvement suggests that this technology is getting closer to prime time, and more teams will begin building experiences around it.

One of the core strengths of video interfaces is their ability to combine multiple sensory channels—visual cues, speech, text overlays, and animations—into a richer, more expressive interaction. Video agents can mimic humanlike expressions and gestures, adding emotional nuance to their communication. For example, an AI-powered customer service avatar might use facial expressions and hand gestures to reassure a frustrated customer, complementing its spoken responses with visual empathy.

However, video interfaces come with technical and design challenges. High-quality video interactions require significant processing power and bandwidth, which can introduce lag or pixelation, undermining the user experience. The uncanny valley remains a risk—if an agent's facial expressions, gestures, or lip-syncing feel slightly off, it can create discomfort rather than engagement. Additionally, privacy concerns are amplified with video agents, as users may feel uneasy about sharing visual data with AI systems.

Looking ahead, video interfaces are poised for significant growth, especially as improvements in rendering, real-time animation, and bandwidth optimization address current limitations. In the near future, expect to see agents embedded seamlessly into virtual meetings, augmented reality (AR) overlays, and digital customer service avatars.

When thoughtfully executed, video interfaces offer an engaging, humanlike dimension to agent interactions, enhancing clarity, emotional connection, and overall effectiveness. As technology advances, video-based agents are set to play a larger role in industries such as telehealth, education, remote collaboration, and interactive entertainment, reshaping how humans and agents communicate in immersive digital spaces.

## Combining Modalities for Seamless Experiences

While each interaction modality—text, graphical interfaces, voice, and video —has its own strengths and limitations, the most compelling agentic experiences often combine multiple modalities into a single, cohesive user journey. Users don't think in terms of modality boundaries; they simply want to achieve their goals as effortlessly and naturally as possible. The ability to move seamlessly across modalities—maintaining state and context throughout —is a hallmark of great agent system design.

For example, a user might begin interacting with an agent via voice while driving, continue the conversation on their phone through text while walking into a meeting, and later review a graphical dashboard summarizing results on their laptop. In another scenario, a voice assistant might read out a summary of an analytics report before emailing a detailed, text-based version with accompanying charts for later reference. This fluid transition between modalities preserves user context, respects situational constraints, and delivers the right interaction style at each moment.

Designing for modality fluidity requires careful state management and context persistence so that information, task progress, and user preferences are never lost in transition. Agents must also adapt their communication style to suit each modality—for example, delivering concise spoken summaries while providing more detailed textual outputs for review.

This is an exciting time for the field of human-computer interaction. Recent advances in foundation models, multimodal architectures, and agent orchestration are unlocking entirely new ways of interacting with intelligent systems. For the first time, it is technically feasible to build agents that engage users across text, voice, images, and video in a single, unified workflow.

However, while the technology frontier is expanding rapidly, it is critical to remember that core UX and product principles remain unchanged. Building

successful agent experiences isn't about showcasing the latest modality integrations or generative UI capabilities for their own sake. It is about understanding users deeply, meeting them where they are, and creating intuitive, trustworthy, and delightful experiences that solve real problems in their lives.

The best products are not those that merely demonstrate technological sophistication, but those that use technology to amplify human capability in elegant and unobtrusive ways. As we continue to push the boundaries of modality design, let us stay grounded in the timeless goal of great product design: creating tools that people love to use, that make their lives easier, and that empower them to achieve what matters most.

## The Autonomy Slider

A critical yet often overlooked dimension in UX design is the level of autonomy granted to agents. As Andrej Karpathy described, effective agentic systems should allow users to smoothly adjust an agent's autonomy—from fully manual control to partial automation to fully autonomous operation. This concept, often called an autonomy slider, empowers users to choose how much control they wish to retain versus delegate at any given time. Figure 3-4 illustrates a simple example of an autonomy slider interface, enabling users to set the agent to "Manual," "Ask," or "Agent" mode depending on their task, trust, and context.
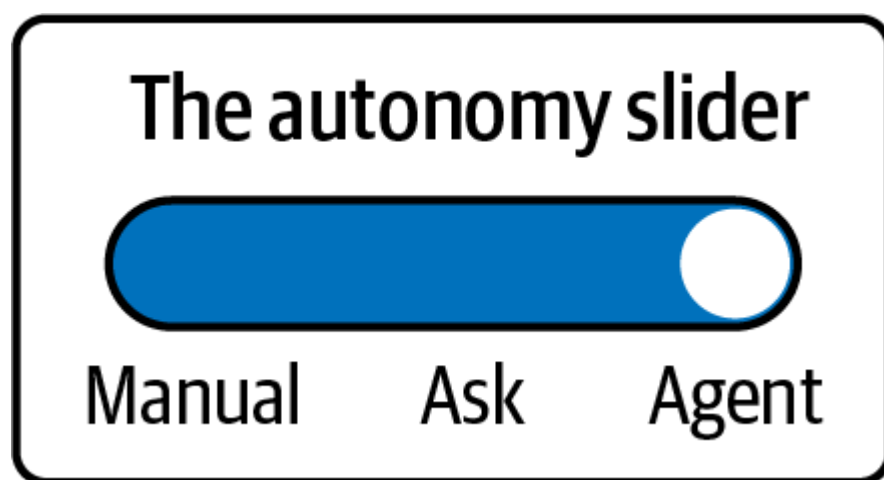


Figure 3-4. The autonomy slider enables users to adjust an agent's level of independence, ranging from fully manual control, to assisted "Ask" mode, to fully autonomous agent execution. This flexibility builds user trust by aligning system behavior with user preferences, task complexity, and context.

Different users, tasks, and contexts demand different degrees of agent autonomy. In some situations, users prefer full manual control to ensure precision, while in others, they may want to offload routine or complex tasks

entirely to the agent. Critically, these preferences are not static; they evolve with user trust, task familiarity, stakes, and workload. For example:

*Manual*

The developer writes all code themselves without agent assistance. The IDE acts purely as an editor with syntax highlighting and linting but no AI-driven suggestions.

*Ask (assisted)*

The agent proactively suggests code completions, refactors, or documentation snippets, but the developer reviews and accepts each suggestion before it is applied. This mode speeds up development while keeping the human fully in control.

*Agent*

The agent autonomously performs certain tasks, such as applying standard refactors, fixing linter errors, or generating boilerplate code files based on project conventions without requiring individual approvals. The developer is notified of changes but does not need to approve each action.

These three modes demonstrate how an autonomy slider empowers developers to balance control and efficiency within a single interface. The same principle applies beyond software development. For example, in a customer support platform:

*Manual*

Human agents handle all incoming customer queries themselves. The AI is inactive or used only for backend analytics, not frontline interactions.

*Ask (assisted)*

The agent drafts suggested replies to customer messages, surfacing recommended responses, policy references, or troubleshooting steps. The human agent reviews, edits if necessary, and approves the reply before sending. This accelerates response time while maintaining human judgment.

*Agent*

The agent autonomously handles routine queries—such as password resets, order tracking, or FAQs—without human intervention, escalating only complex or sensitive issues to human agents. Users are notified of agent actions but do not need to approve each message for standard interactions.

These three modes coexist within the same customer support system, empowering teams to adjust autonomy based on query complexity, customer profile, and organizational trust in AI. This same autonomy slider pattern can extend to any field where workflows benefit from fluidly shifting between manual execution, AI assistance, and full agentic automation. This spectrum of autonomy must be consciously designed into agent experiences. Without it, agents risk feeling either underpowered (if they require too much manual input) or overbearing (if they act without user consent in sensitive contexts). To integrate an autonomy slider effectively, consider the following design principles:

*Expose degrees of autonomy clearly*

Users should understand the available levels of agent independence, from manual to assisted to autonomous. Label these modes in intuitive language, such as "Manual," "Assist," and "Auto," and explain their implications.

*Enable seamless transitions*

Users must be able to shift between autonomy levels effortlessly as their confidence, context, or workload changes. For instance, a toggle or slider in the interface should offer a quick transition from review mode to auto-approve mode.

*Provide predictable and transparent behavior at each level*

Each autonomy level should have well-defined behaviors. In partial automation, for example, the agent may draft an output but require explicit user approval before execution. In full autonomy, it should still provide status updates and options to intervene.

*Communicate the risks and benefits of each level*

Users should be aware of what they gain or risk by increasing agent autonomy. For critical tasks, it may be advisable to require an explicit user confirmation before enabling full autonomy.

*Adapt autonomy based on user trust and competence*

Intelligent systems can gradually suggest higher autonomy levels as users gain trust and as the agent demonstrates reliability. For example, after 10 successful uses in manual mode, the system might suggest trying assist mode to save time.

Importantly, the autonomy slider is not merely a feature—it is a trust-building mechanism. By giving users control over how much autonomy an agent exercises, systems communicate respect for user expertise and agency. It avoids the common pitfall of "one-size-fits-all" autonomy that either overwhelms or underutilizes user potential. Always ask: how easily can my users move between manual, assisted, and fully autonomous modes? The answer to this question will shape whether your agent is adopted as a reliable partner or sidelined as an untrusted tool.

# Synchronous Versus Asynchronous Agent Experiences

Agent systems can operate in synchronous or asynchronous modes, each offering distinct advantages and challenges. In synchronous experiences, interactions occur in real time, with immediate back-and-forth exchanges between the user and the agent. These experiences are common in chat interfaces, voice conversations, and real-time collaboration tools, where quick responses are essential for maintaining flow and engagement. In contrast, asynchronous experiences enable agents and users to operate independently, with communication occurring intermittently over time. Examples include email-like interactions, task notifications, or agent-generated reports delivered after a process has completed.

The choice between synchronous and asynchronous designs depends heavily on the nature of the task, user expectations, and operational context. While synchronous agents excel in tasks requiring instant feedback or live decision making, asynchronous agents are better suited for workflows where tasks may take longer, require background processing, or don't demand the user's

constant attention. Striking the right balance between these modes—and managing when agents proactively engage users—can greatly influence user satisfaction and the overall effectiveness of the system. Both are useful and valid patterns, but it is highly recommended to choose which experiences fall into which category, so that users do not end up waiting for a pinwheel to spin.

## Design Principles for Synchronous Experiences

Synchronous agent experiences thrive on immediacy, clarity, and responsiveness. Users expect agents in these settings to respond quickly and maintain conversation flow and context without noticeable delays. Whether in a live chat, voice call, or real-time data dashboard, synchronous interactions demand low latency and context awareness to avoid frustrating pauses or repetitive questions.

Agents in synchronous environments should prioritize clarity and brevity in their responses. Long-winded explanations or overly complex outputs can break the rhythm of real-time interactions. Additionally, turn-taking mechanics—knowing when to respond, when to wait, and when to escalate—are critical for maintaining a natural and productive conversation flow. Visual cues, like typing indicators or progress spinners, can reassure users that the agent is actively processing their input.

Error handling is equally important in synchronous designs. Agents must gracefully recover from misunderstandings or failures without derailing the interaction. When uncertainty arises, synchronous agents should ask clarifying questions or gently redirect users rather than making risky assumptions. These principles create a smooth, intuitive experience that keeps users engaged and maintains context without unnecessary friction.

## Design Principles for Asynchronous Experiences

Asynchronous agent experiences prioritize flexibility, persistence, and clarity over time. These interactions often occur in contexts where immediate responses aren't necessary, such as when agents are processing long-running tasks, preparing detailed reports, or monitoring background events.

Effective asynchronous agents must excel at clear communication of task status and outcomes. Users should always understand what the agent is doing,

what stage a task is in, and when they can expect an update. Notifications, summaries, and well-structured reports become key tools for maintaining transparency. For example, an agent generating an analytical report might notify the user when processing begins, provide an estimated completion time, and deliver a concise, actionable summary when finished.

Context management is another critical design principle for both asynchronous and synchronous agents. Because there may be long delays between user-agent interactions, agents must retain and reference historical context seamlessly. Users shouldn't need to repeat information or retrace previous steps when returning to an ongoing task. We'll cover this in more detail in [Chapter 6](#) on memory.

Lastly, asynchronous agents must manage user expectations effectively. Clear timelines, progress indicators, and follow-up notifications prevent frustration caused by uncertainty or lack of visibility into an agent's work.

## Finding the Balance Between Proactive and Intrusive Agent Behavior

One of the most delicate aspects of agent design—whether synchronous or asynchronous—is determining when and how agents should proactively engage users. Proactivity can be immensely helpful, such as when an agent alerts a user to an urgent issue, suggests an optimization, or provides a timely reminder. However, poorly timed notifications or intrusive behaviors can frustrate users, disrupt their workflow, or even cause them to disengage entirely.

The key to balancing proactivity lies in context awareness and user control. Agents should understand the user's current focus, level of urgency, and communication preferences. For instance, a proactive alert during a high-stakes video meeting might be more disruptive than helpful, while a notification about a completed task delivered via email might be perfectly appropriate.

Agents should also prioritize relevance when proactively reaching out. Notifications and suggestions must add genuine value—solving problems or providing insights rather than adding noise. Additionally, users should have

control over notification frequency, channels, and escalation thresholds, enabling them to customize agent behavior to suit their needs.

Striking this balance isn't just about technical capability—it's about empathy for the user's workflow and mental state. Well-designed agents seamlessly weave proactive engagement into their interactions, enhancing productivity and reducing friction without becoming overbearing.

# Context Retention and Continuity

Ensuring context retention and continuity across user interactions is an important aspect of designing effective agent systems. Whether an agent is guiding a user through a multistep workflow, continuing a paused conversation, or adjusting its behavior based on past interactions, its ability to maintain context directly impacts usability, efficiency, and user trust.

While context retention is a technical capability, it is fundamentally a UX consideration because it determines whether users experience the agent as a cohesive, attentive collaborator or as a disconnected tool that forces them to repeat themselves. From the user's perspective, memory creates a sense of continuity, personalization, and intelligence. If an agent remembers previous interactions, user preferences, or in-progress tasks, it can seamlessly continue conversations and workflows, reducing cognitive load and frustration.

Implementation approaches directly shape UX. A purely client-side context (e.g., stored in browser memory) may feel fast within a session but loses continuity across devices or logins, undermining seamless UX. A purely server-side context (e.g., stored in a database tied to user ID) enables long-term memory and cross-device experiences but can introduce latency or privacy considerations. A hybrid approach—maintaining short-term context on the client side for responsiveness and persisting long-term context on the server side for continuity—often achieves the best UX balance. Choosing the right strategy depends on the user journey, privacy requirements, and level of personalization intended. Ultimately, context is UX: it is how an agent remembers, adapts, and responds in ways that make it feel human-centered and supportive rather than stateless or mechanical.

Effective context retention requires agents to manage both short-term and long-term memory effectively. Short-term memory enables an agent to hold

details within an ongoing session, such as remembering the specifics of a question or instructions given moments earlier. Long-term memory, on the other hand, enables agents to retain preferences, past interactions, and broader user patterns across multiple sessions, enabling them to adapt over time.

However, context management introduces challenges. Data persistence, privacy concerns, and memory limitations must all be carefully addressed. If an agent loses track of context mid-task, the user experience can feel disjointed, repetitive, and frustrating. Conversely, if an agent retains too much context or stores unnecessary details, it risks becoming unwieldy or even breaching user privacy.

In the next section, we'll explore two key facets of context retention and continuity: maintaining state across interactions, and personalization and adaptability—both essential for delivering fluid, intuitive, and user-centric agent experiences.

## Maintaining State Across Interactions

State management is the foundation of context continuity in agent systems. For an interaction to feel seamless, an agent must accurately track what has happened so far, what the user intends to achieve, and what the next logical step is. This is particularly important in multiturn conversations, task handoffs, and workflows with intermediate states, where losing context can result in frustration, inefficiency, and abandonment of tasks.

Effective state management depends on how the system identifies and tracks users or sessions. For logged-in users, state can be tied directly to their user accounts, enabling memory persistence across devices and sessions. For anonymous interactions, maintaining context typically requires a session identifier—such as a cookie or token—to track the conversation between the client and server.

As agent systems scale to thousands or millions of users, session state should not reside only in memory. Persisting state in a database or distributed cache ensures continuity across server restarts, enables load balancing, and supports multidevice experiences. The choice between user-based memory (persistent, personalized) and session-based memory (ephemeral, session-scoped) depends on your application's privacy requirements, user expectations, and operational architecture. Regardless of implementation, robust identification

and storage strategies are fundamental to delivering seamless, context-aware agent experiences at scale.

Agents can maintain state through short-term session memory, where details of the ongoing interaction—such as a user's recent commands or incomplete tasks—are temporarily stored until the session ends. In more advanced systems, persistent state management enables agents to resume tasks across multiple sessions so that users can pick up where they left off, even after hours or days have passed.

Effective state retention requires clear session boundaries, data validation, and fallback mechanisms. If an agent forgets context, it should gracefully recover by asking clarifying questions rather than making incorrect assumptions. Additionally, state data must be managed securely and responsibly, especially when it involves sensitive or personally identifiable information.

When done well, maintaining state enables agents to guide users through complex tasks without unnecessary repetition, reduce cognitive load, and create a sense of ongoing collaboration. Whether an agent is helping a user book travel accommodations, troubleshoot a technical issue, or manage a multistep approval process, effective state management ensures interactions remain smooth, logical, and productive.

## Personalization and Adaptability

Personalization goes beyond merely remembering context—it involves using past interactions and preferences to tailor the agent's behavior, responses, and recommendations to individual users. An adaptable agent doesn't just maintain state; it learns from previous exchanges to deliver increasingly refined and relevant outcomes. Personalization can take multiple forms:

*Preference retention*

> Remembering user settings, such as notification preferences or commonly chosen options

*Behavioral adaptation*

> Adjusting response style or interaction flow based on observed user patterns

*Proactive assistance*

Anticipating user needs and offering suggestions based on past behavior

For example, an agent assisting with project management might recognize a user's preferred task-tracking style and adapt its notifications or summaries accordingly. Similarly, a customer service agent might adjust its tone and verbosity based on whether the user prefers concise answers or detailed explanations.

However, personalization comes with challenges. Privacy concerns must be carefully managed, with transparent communication about what data is being stored and how it is being used. Additionally, agents must strike a balance between being helpfully adaptive and overly persistent—users should always have the option to reset or override personalized settings.

The best personalization feels invisible yet impactful, where the agent subtly improves the user experience without drawing attention to its adjustments. At its peak, personalization creates an experience where users feel understood and supported, as if the agent is a thoughtful collaborator rather than a mechanical tool.

# Communicating Agent Capabilities

One of the most critical aspects of designing effective agent experiences is ensuring users understand what the agent can do and how to interact with it effectively. While backend agent design determines what functions an agent supports, the user experience determines whether those capabilities are discoverable, intuitive, and usable in practice. In traditional applications, discoverability is straightforward: menus, buttons, and interface elements visually communicate available actions. In agentic systems, especially those using text or voice interfaces, the absence of visible affordances often leaves users guessing what the agent can and cannot do.

Effective agent UX addresses this challenge by proactively communicating capabilities through the interface itself. For example, many chat-based agents include suggested action buttons below the input field, highlighting common or contextually relevant actions such as "Track order," "Generate summary," or "Create meeting note." These buttons serve as visual affordances, guiding

users toward supported workflows without requiring them to remember specific commands or guess what is possible. Similarly, onboarding tutorials or first-use walkthroughs can introduce users to an agent's core functions, helping them build confidence early on.

Another useful pattern is the inclusion of expandable menus or capability cards that list available functions in a structured way. In a graphical agent interface, for instance, a sidebar might contain sections for data retrieval, analysis, summarization, and workflow automation. This mirrors the menu structures that users expect in traditional apps while communicating the breadth of agent capabilities upfront. Dynamic suggestions, where the system recommends actions based on user input, also help bridge the gap between open-ended natural language and structured tool invocation. If a user begins typing "book…," the agent might suggest "Book meeting with [name]," "Book conference room," or "Book travel," anticipating intent and making actions easier to execute.

In systems relying primarily on open-ended text input, agents themselves must communicate their capabilities clearly in conversation. This can include proactive introductions when a session begins, such as: "Hi, I can help you generate content, analyze data, or summarize documents. What would you like to do today?" When users request actions beyond current capabilities, the agent should not simply reject the request but provide alternatives: "I can't process payments directly, but I can update your billing preferences or connect you with an agent who can assist." Such responses reduce user frustration while reinforcing the agent's utility.

While it is important to surface capabilities, it is equally critical not to overwhelm users with too many options at once. Effective designs prioritize progressive disclosure, showing core capabilities initially and revealing advanced features as users become more comfortable. Contextual relevance also plays a key role. Displaying the most likely actions based on current user inputs, historical behavior, or workflow stage ensures the agent feels supportive rather than cluttered. Visual grouping and clear hierarchy within menus or suggested actions help users navigate available options efficiently.

These principles apply across modalities. In text-based chat interfaces, quick-reply buttons and example prompts improve clarity. In graphical dashboards, capability menus and tooltips communicate functions without crowding the interface. Voice agents must balance brevity with clarity, listing only a few

high-priority options at a time to avoid cognitive overload. Generative UI systems can combine natural language and dynamically generated visual outputs to make available capabilities immediately visible and actionable.

Ultimately, communicating agent capabilities is not merely about stating what the agent can do; it is about designing an experience that empowers users to harness those capabilities confidently and efficiently. When users understand an agent's scope and limitations, they are far more likely to engage productively, trust its outputs, and integrate it into their workflows. Thoughtful UX design turns invisible functions into visible affordances, transforming agents from opaque black boxes into transparent, collaborative digital partners.

## Communicating Confidence and Uncertainty

Agents often operate in probabilistic environments, generating outputs based on statistical models rather than deterministic rules. As a result, not every response or action carries the same degree of confidence. Communicating uncertainty effectively is essential for building user trust and helping users make informed decisions.

Confidence levels can be expressed in several ways:

*Explicit statements*

"I'm 90% certain this is the correct answer."

*Visual cues*

Icons, color-coded alerts, or confidence meters in graphical interfaces.

*Behavioral adjustments*

Offering suggestions rather than firm recommendations when confidence is low.

Agents must avoid appearing overly confident when uncertainty is high—users are quick to lose trust if an agent confidently delivers an incorrect or misleading response. Similarly, excessive hedging in low-stakes interactions can make an agent appear hesitant or unreliable.

Communicating confidence and uncertainty isn't just about sharing probabilities; it's about framing responses in a way that aligns with user expectations and the stakes of the interaction. In critical contexts, transparency is nonnegotiable, while in low-stakes settings, confidence can be presented more casually.

## Asking for Guidance and Input from Users

No agent, no matter how advanced, can perfectly interpret ambiguous, vague, or conflicting user inputs. Instead of making risky assumptions, agents must know when to ask clarifying questions or seek user guidance. This ability transforms potential errors into opportunities for collaboration.

Effective agents are designed to ask focused, helpful questions when they encounter ambiguity. For example, if a user says "Book me a ticket to Chicago," the agent might respond with "Would you like a one-way or round-trip ticket, and do you have preferred travel dates?" Instead of defaulting to a generic response or making incorrect assumptions, the agent uses the opportunity to refine its understanding.

The way agents ask for guidance also matters. Questions should be clear, polite, and context-aware, avoiding robotic or repetitive phrasing. If the user has already answered part of the question earlier in the conversation, the agent should reference that context rather than starting from scratch.

Additionally, agents should be transparent about why they're asking for clarification. A simple explanation, like "I need a bit more information to proceed accurately," helps users understand the rationale behind the question.

Finally, agents should avoid asking too many questions at once—this can overwhelm users and make the interaction feel like an interrogation. Instead, they should sequence questions logically, addressing the most critical ambiguities first.

When agents confidently ask for guidance and input, they transform uncertainty into productive collaboration, empowering users to guide the agent toward successful outcomes while maintaining a sense of partnership and shared control.

# Failing Gracefully

Failure is inevitable in agentic systems. Whether due to incomplete data, ambiguous user input, technical limitations, or unexpected edge cases, agents will encounter scenarios where they cannot fulfill a request or complete a task. However, how an agent handles failure is just as important as how it handles success. A well-designed agent doesn't just fail—it fails gracefully, minimizing user frustration, preserving trust, and providing a clear path forward.

At its core, graceful failure involves acknowledging the issue transparently, offering a helpful explanation, and suggesting actionable next steps. For instance, if an agent cannot find an answer to a query, it might respond with "I couldn't find the information you're looking for; would you like me to escalate this to a human representative?" instead of producing an incorrect or nonsensical response.

Agents should also be designed to anticipate common points of failure and have predefined fallback mechanisms in place. For example, if a voice-based agent struggles to understand repeated user inputs, it might switch to a text-based option or provide a clear explanation, such as: "I'm having trouble understanding your request. Could you please try rephrasing it or typing your question instead?"

In multistep tasks, state preservation is equally important when an agent encounters failure. Instead of requiring the user to restart from scratch, the agent should retain progress and allow the user to pick up where they left off once the issue is resolved. This prevents unnecessary repetition and frustration.

Another critical aspect of graceful failure is apologetic and empathetic language. When something goes wrong, the agent should acknowledge the failure in a way that feels human and considerate, avoiding cold or overly technical error messages. For example: "I'm sorry; something went wrong while processing your request. Let me try again or connect you with someone who can help."

Additionally, agents should provide clear paths to resolution. Whether it's offering troubleshooting steps, escalating to a human operator, or directing the

user to an alternative resource, users should always know what options are available to them when the agent encounters a roadblock.

Lastly, agents must learn from their failures whenever possible. Logging failure points, analyzing recurring issues, and feeding these insights back into the development process can help reduce the frequency of similar failures in the future. Agents that improve iteratively based on their failure patterns will become increasingly resilient and reliable over time.

In summary, failing gracefully is about maintaining user trust and minimizing frustration even when things don't go as planned. By being transparent, empathetic, and action-oriented, agents can turn failures into opportunities to strengthen their relationship with users, demonstrating reliability even in moments of imperfection.

## Trust in Interaction Design

Trust is gained in drops and lost in buckets. This certainly applies to agentic systems as well. Without it, even the most advanced agent systems will struggle to gain user acceptance, regardless of their capabilities. Transparency and predictability are two of the most powerful tools for building and maintaining trust between agents and users. Users need to understand what an agent can do, why it made a particular decision, and what its limitations are. This clarity fosters confidence, reduces anxiety, and encourages productive collaboration.

Transparency begins with clear communication of agent capabilities and constraints. Users should never have to guess whether an agent can handle a task or if it is operating within its intended scope. When agents provide explanations for their actions—whether it's how they arrived at a recommendation, why they declined a request, or how they interpreted an ambiguous instruction—they give users visibility into their reasoning. This isn't just about building trust; it also helps users refine their instructions, improving the quality of future interactions.

Predictability complements transparency by ensuring that agents behave consistently across different scenarios. Users should be able to anticipate how an agent will respond based on prior interactions. Erratic or inconsistent behavior, even if technically correct, can quickly erode trust. For example, if

an agent suggests a cautious approach in one context but appears overly confident in a nearly identical scenario, users may start to question the agent's reliability.

However, transparency does not mean overwhelming the user with unnecessary details. Users don't need to see every step of the agent's reasoning process—they just need enough insight to feel confident in its actions. Striking this balance requires thoughtful interface design, using visual cues, status messages, and brief explanations to communicate what's happening without causing cognitive overload.

When trust and transparency are prioritized, agent systems become more than just tools—they become reliable collaborators. Users feel confident delegating tasks, following agent recommendations, and relying on their outputs in both casual and high-stakes scenarios. In the remainder of this section, we'll explore two key components of trust-building: ensuring predictability and reliability in agent behavior.

Predictability and reliability are foundational to trust. Users must be able to count on agents to behave consistently, respond appropriately, and handle errors gracefully. Agents that act erratically, give conflicting outputs, or produce unexpected behavior—even if occasionally correct—can quickly undermine user confidence.

Reliability begins with consistency in agent outputs. If a user asks an agent the same question under the same conditions, they should receive the same response. In cases where variability is unavoidable (e.g., probabilistic outputs from language models), agents should clearly signal when an answer is uncertain or context-dependent.

Agents must also handle edge cases thoughtfully. For example, when they encounter incomplete data, conflicting instructions, or ambiguous user input, they should respond predictably—either by asking clarifying questions, providing a neutral fallback response, or escalating the issue appropriately.

Another critical aspect of reliability is system resilience. Agents should be designed to recover from errors, maintain state across interruptions, and prevent cascading failures. For example, if an agent loses connection to an external API, it should notify the user, explain the issue, and offer a sensible next step rather than silently failing or producing misleading outputs.

Lastly, reliability is about setting and meeting expectations consistently. If an agent claims it can handle a specific task, it must deliver on that promise every time. Misaligned expectations—where agents overpromise and underdeliver—can cause more damage to user trust than simply admitting limitations up front.

When agents behave predictably and reliably, they become dependable digital partners, empowering users to trust their outputs, delegate tasks confidently, and rely on them for critical decisions.

# Conclusion

Designing exceptional user experiences for agent systems goes far beyond technical functionality—it requires an understanding of how humans interact with technology across different modalities, contexts, and workflows. Whether through text, graphical interfaces, voice, or video, each interaction modality carries its own strengths, trade-offs, and unique design considerations. Successful agent experiences are those where the modality aligns seamlessly with the user's task, environment, and expectations.

Synchronous and asynchronous agent experiences present distinct design challenges, requiring thoughtful approaches to timing, responsiveness, and clarity. Synchronous interactions demand immediacy and conversational flow, while asynchronous interactions excel in persistence, transparency, and thoughtful notifications. Striking the right balance between proactive assistance and intrusive interruptions remains one of the most delicate aspects of agent design.

Exceptional agents seamlessly retain context and adapt to users, remembering critical details across interactions and adapting intelligently to user preferences. This ability not only reduces cognitive load but also fosters a sense of continuity and collaboration, transforming agents from isolated tools into reliable digital partners. Some common patterns to keep in mind:

*Communicate capabilities clearly*

Show users what the agent can do through onboarding, suggestions, or buttons.

*Combine modalities thoughtfully*

Align text, GUI, voice, or video with the task and user context.

*Retain context thoughtfully*

Maintain relevant conversation state without overwhelming memory or violating privacy.

*Handle errors gracefully*

Provide clear, polite fallbacks when the agent can't fulfill a request.

*Build trust*

Be transparent about limitations, confidence, and reasoning.

Equally important is how agents communicate their capabilities, limitations, and uncertainties. Clear expectations, honest confidence signals, and thoughtful clarification questions create trust, reduce frustration, and prevent misunderstandings. Agents must also know how to fail gracefully, guiding users toward alternative solutions without leaving them stranded or confused.

Finally, building trust through predictability, transparency, and responsible design choices ensures that users can rely on agents. Trust is earned not just through success but also through how agents handle ambiguity, failure, and recovery.

As the agent landscape continues to shift and expand, designers and developers must remain agile—continually reevaluating interaction paradigms, adapting to new multimodal capabilities, and experimenting with novel UX patterns. The design patterns described here provide a robust starting point, but the future of agentic UX will be shaped by rapid innovation in modalities, context management, and human-agent collaboration. In the years ahead, agent systems will continue to evolve, becoming more deeply embedded in our personal and professional lives. The principles outlined in this chapter—focused on clarity, adaptability, transparency, and trust—provide a blueprint for creating agent experiences that are not just functional, but intuitive, engaging, and deeply aligned with human needs.

By prioritizing UX at every stage of development, we can ensure that agents become not just tools, but indispensable partners in our increasingly intelligent digital ecosystems. In Chapter 4, we'll cover tool use, which is how we move from ordinary chatbots to systems that can do real work for users.