# Novice-AI Music Co-Creation via AI-Steering Tools for Deep Generative Models

**Ryan Louie**
Northwestern University
ryanlouie@u.northwestern.edu

**Andy Coenen**
Google PAIR/Big Picture
andycoenen@google.com

**Cheng-Zhi Anna Huang**
Google Magenta
chengzhiannahuang@gmail.com

**Michael Terry**
Google PAIR/Big Picture
michaelterry@google.com

**Carrie Cai**
Google PAIR/Big Picture
cjcai@google.com

## ABSTRACT

While generative deep neural networks (DNNs) have demonstrated their capacity for creating novel musical compositions, less attention has been paid to the challenges and potential of co-creating with these musical AIs, especially for novices. In a needfinding study with a widely used, interactive musical AI, we found that the AI can overwhelm users with the amount of musical content it generates, and frustrate them with its non-deterministic output. To better match co-creation needs, we developed AI-steering tools, consisting of Voice Lanes that restrict content generation to particular voices; Example-Based Sliders to control the similarity of generated content to an existing example; Semantic Sliders to nudge music generation in high-level directions (e.g., happy/sad); and Multiple Alternatives to generate multiple possibilities to choose from. In a summative study (N=21), we discovered the tools not only increased users' trust, control, comprehension, and sense of collaboration with the AI, but also contributed to a greater sense of self-efficacy and ownership of the composition relative to the AI.

## INTRODUCTION

Rapid advances in deep learning have made it possible for artificial intelligence (AI) to actively collaborate with humans to co-create new content [32, 5, 14, 29, 18, 27]. One promising application of machine learning in this space has been the use of generative deep neural network (DNN)-backed systems for creative activities such as poetry writing, drawing, and music creation—experiences that bear intrinsic value for people, but that often require specialized skill sets. For example, by completing a drawing that a user has started [32, 5, 28, 10] or filling in a missing section of a song [23, 19], generative models could enable untrained lay users to take part in creative experiences that would otherwise be difficult to achieve without additional training or specialization [25, 8, 15]. In

this paper, we focus on the needs of music novices co-creating music with a generative DNN model.

While substantial work has focused on improving the algorithmic performance of generative music models, little work has examined what interaction capabilities users actually need when co-creating with generative AI, and how those capabilities might affect the music co-creation experience. Recent generative music models have made it conceivable for novices to create an entire musical composition from scratch, in partnership with a generative model. For example, the widely available Bach Doodle [25] sought to enable anyone on the web to create a four-part chorale in the style of J.S. Bach by writing only a few notes, allowing an AI to fill in the rest. While this app makes it conceivable for even novices with no composition training to create music, it is not clear how people perceive and engage in co-creation activities like these, or what types of capabilities they might find useful.

In a study we conducted to understand the human-AI co-creation process, we found that AI music models can sometimes be quite challenging to co-create with. Paradoxically, the very capabilities that enable such sophisticated models to rival human performance can impede human partnership: users struggled to evaluate and edit the generated music because the system created *too much* content at once; in essence, they experienced *information overload*. They also struggled with the system's *non-deterministic output*: while the output would typically be coherent, it would not always align with users' musical goals at the moment. These findings raise critical questions about how to co-create with an AI that already matches or supercedes human generative capabilities: what user interfaces and interactive controls are important, and what interactive capabilities should be exposed by deep generative neural nets to benefit co-creation?

In this work, we examine what novices need when co-creating music with a deep generative model, then propose and evaluate **AI-steering tools** that enable novice users to iteratively direct the creation process in real-time. For the purposes of this work, we define *novices* as people who have played a musical instrument, but who have little or no experience formally composing music. To ground this research, we developed Cococo (collaborative co-creation), a music editor web-interface for novice-AI co-creation that augments standard generative

music interfaces with a set of AI-steering tools: 1) *Voice Lanes* that allow users to progressively define when and where the AI generates music, before any music is created, 2) an *Example-based Slider* for expressing that the AI-generated music should be more or less like an existing example of music, 3) *Semantic Sliders* that users can adjust to direct the music toward high-level directions (e.g. happier / sadder, or more conventional / more surprising), and 4) *Multiple Alternatives* for the user to select between a variety of AI-generated options. To implement the sliders, we developed a *soft priors* approach that enables the AI to both adhere to surrounding musical context and match user-desired qualities, without needing to retrain the model.

In a summative evaluation with 21 music novices, we found that AI-steering tools not only increased users' trust, control, comprehension, and sense of collaboration with the AI, but also contributed to a greater sense of self-efficacy and ownership of the composition relative to the AI. Beyond improving user attitudes towards the AI, the tools also enabled new user strategies for music co-creation: participants used the tools to divide the music into semantically meaningful components; learn and discover musical structure; debug the music and the AI; and explore the limits of the AI.

In sum, this paper makes the following contributions:

- We discover key needs of music novices when co-creating with a typical generative-DNN music interface, including issues related to AI-induced *information overload* and its *non-deterministic output*.

- We present the design and implementation of AI-steering tools that enable users to progressively guide the co-creation process in real-time, employing *soft priors* to influence the AI's content generation.

- We find in a summative study with 21 users that the tools increase users' sense of ownership of the composition relative to the AI, while increasing trust, controllability, and comprehensibility of the AI.

- We describe new user strategies for co-creating with AI using these tools, such as developing new insights into composition strategies, isolating the cause of musical glitches, and exploring the limits of the AI. We also uncover novice considerations of agency and collaboration when co-creating with AI.

Taken together, these findings inform the design of future human-AI interfaces for co-creation.

## RELATED WORK

### Human-AI Co-creation
The acceleration of AI capabilities has renewed interest in how AI can enable human-AI co-creation in domains such as drawing [32, 5, 28, 10], creative writing [14, 4], design ideation [29], video game content generation [18], and dance [27]. For example, an AI might flesh out a half-sketched drawing [32], write the next paragraph of a story [4], or add an image to a design mood board [29]. Across this range of prior work, a core challenge has been developing collaborative

AI agents that can adapt their actions based on the goals and behaviors of the user. To this end, some systems design the AI to generate output conditioned upon the surrounding context of human-generated content [10, 4, 14], while others leverage user feedback to better align AI behavior to user intents [18, 29, 5]. Research has also observed that users desire to take initiative in their partnership with AI [32], with controllability and comprehensibility being key challenges to realizing this vision [1]. Building on this need, our work enables users to express their preferences to an AI collaborator through a variety of means.

Much of the prior work in this space has focused on the domains of drawing or writing. There have been relatively fewer HCI efforts examining generative DNN music agents of similar prowess. Building on prior work examining AI as a peer in the creative process, our work contributes to the broader literature by investigating human-AI co-creation in music.

### Interactive Interfaces for ML Music Models
To support music makers in the composition process, researchers have developed ML-powered interfaces and devices that map user inputs to musical structures so users can interactively explore musical variations. Examples of such systems include those that allow users to find chords to accompany a melody [37, 17], experiment with adventurous chord progressions [24, 13], use custom gestural inputs to interpolate between synthesizer sounds [12], or turn free-hand sketches into harmonious musical textures [11].

More recently, progress in generative DNNs has introduced fully-generative music interfaces capable of performing auto-completion given a seed of user-specified notes [19, 25, 35]. Beyond supporting single sub-components, these systems can produce full scores that automatically mesh well with local and distant regions of music. Thus, there is potential to now support users in a wide range of musical tasks (e.g., harmonizing melodies, elaborating existing music, composing from scratch), all within one interface. While recent research has made these fully-generative interfaces increasingly available to musicians and novices alike [19, 35, 25, 7], there has been relatively little HCI work examining how to design interactions with these contemporary models to ensure they are effective for co-creation, especially for novices. Our research contributes an integrative understanding of how interfaces to these capable AIs can be designed and used, how these capabilities affect the composing experience, and users' attitudes towards AI co-creation.

### Deep Generative Music Models
As their name implies, generative deep neural networks can synthesize content. Research has demonstrated the potential for modeling and synthesizing music, ranging from single-voice sequences [9] and multi-part music [15, 30], to music with variable parts at each time step [2] and music with long-term structure over minutes [26, 33, 21].

In contrast to models that (typically) generate music chronologically from left to right, *in-filling* models can more flexibly support co-creation by allowing users to specify regions at any point in the music, then auto-filling those gaps. Examples

include DeepBach [19] and Coconet [23], both trained on four-part Bach Chorales. Researchers have also created models designed to support interaction mechanisms that grant users more control. For example, there are emerging approaches aimed at learning a continuous latent space so that users can interpolate between music [34], or explore a space of musical alternatives [6]. In our work, we adopt *soft priors* (described more fully below) as a general approach that provides additional ways for users to direct their exploration. In contrast to hard constraints, our approach allows DNNs to simultaneously consider the original context and the newly added soft priors, without needing to re-train the model.

## FORMATIVE NEEDFINDING STUDY

Our research focuses on enabling music composition novices to engage more creatively with music, without the prerequisite understanding of musical theory and composition. Thus, we conducted a 45 minute formative interview and elicitation study with 11 music novices to understand (1) their motivations and needs for creating music themselves and (2) challenges in co-creating with AI composing tools. We recruited participants from our institution using mailing lists and word-of-mouth, screening for individuals who had played a musical instrument at some point in their life: 9 participants had five or more years of experience playing a musical instrument; 8 had no formal experience in composition and had informally experimented with musical arrangements using music software or improvising on an instrument; and 2 had tried creating a small composition as part of a music theory class assignment.

### Motivations and Needs for Creating Music

Our participants reported the desire to create music to complement or enrich existing personal artifacts or experiences, such as creating an accompaniment to a short personal video or photo album, a composition inspired by a poem, or a theme-song for a friend or loved one. Participants who had attempted creating music on their own encountered challenges due to their lack of training in music theory and composition. Oftentimes, they knew something needed to be created or fixed (e.g., adding harmonies), but lacked the expertise to identify the issue, a strategy for solving the problem, and/or the ability to generate viable solutions. Each of these challenges suggest specific ways AIs could aid users and make them more capable.

### Challenges in Co-Creating with Generative DNNs

In the second half of the study, we conducted an elicitation to understand challenges when interacting with a deep generative model to compose music. The interface mirrored the generative *infilling* capabilities found in contemporary interfaces for deep generative models [25], where users can manually draw notes and request the AI to fill in the remaining voices and measures, or erase any part of the music and request the AI to fill in the gap. Overall, we found that users struggled to evaluate the generated music and express desired musical elements, due to *information overload* and *non-deterministic output*.

### Information Overload

While the deep generative models were capable of infilling much of the song based on only a few notes from the user, participants found the amount of generated content overwhelming to unpack, evaluate, and edit. Specifically, they had difficulty determining why a composition was off, and expressed frustration at the inability to work on smaller, semantically meaningful parts of the composition. For example, one user struggled to identify which note was causing a discordant sound after multiple generated voices were added to their original: *"It was difficult because all the notes were put on the screen already... I can identify places where it doesn't sound very good, but it's actually hard to identify the specific note that is off."* Some participants naturally wanted to work on the composition *"bar-by-bar or part-by-part"*; in contrast to expectations, the generated output felt like it *"skipped a couple steps"* and made it difficult to follow all at once: *"Instead of giving me four parts of harmony, can it just harmonize one? I can't manage all four at once."*

### Non-deterministic output

Even though the AI was capable of generating notes that were technically coherent to the context of surrounding notes provided by users, the stochastic nature of the system meant that its output did not always match the user's current musical objectives. For example, a participant who had manually created a dark, suspenseful motif was dismayed with how the generated notes were misaligned with the original feeling of the motif: *"the piece lost the essence of what I was going for. While it sounds like nice music to play at an upscale restaurant, the sense of climax is not there anymore."* Even though what was produced sounded harmonious to the user, they felt incapable of giving feedback about their goal in order to constrain the kinds of notes the model generated. Despite being *technically* aligned to context, the music was *musically* mis-aligned with user goals. As a result, participants wished there were ways to go beyond randomly *"rolling dice"* to generate a desired sound, and instead control the generation based on relevant musical objectives.

## COCOCO

Based on identified user needs, we developed Cococo (collaborative co-creation), a music editor web-interface for novice-AI co-creation that augments standard generative music interfaces with a set of *AI steering tools* (Figure 1). Cococo builds on top of Coconet [22], a state-of-the-art deep generative model trained on 4 part harmony that accepts incomplete music as input and outputs complete music. Coconet works with music that can have 4 parts or *voices* playing at the same time (represented by **S**oprano **A**lto **T**enor **B**ass), are 2-measures long or 32 *timesteps* of sixteenth-note beats, and where each voice can take on any one of 46 *pitches*. Coconet is able to *infill* any section of music, including gaps in the middle or start of the piece. To mirror the most recent interfaces backed by these infill capabilities [7, 19], Cococo contains an *infill mask* feature, with which users can crop a passage of notes to be erased using a rectangular mask, and automatically infill that section using AI. Users can also manually draw and edit notes.
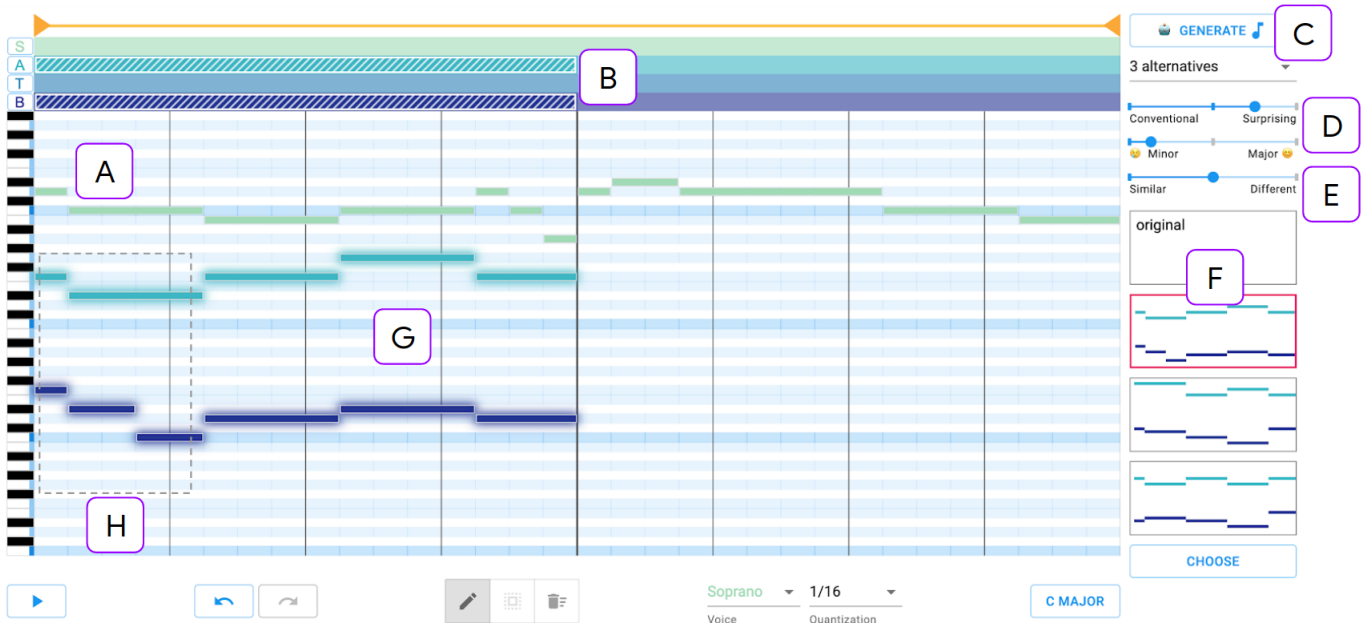
**Figure 1. Key components of Cococo: users can manually write some notes (A), specify which voices and in which time range to request AI-generated music using Voice Lanes (B), click Generate (C) to ask the AI to fill in music given the existing notes on the page, use Semantic Sliders (D) to steer or adjust the AI's output along semantic dimensions of interest (e.g. more surprising), use the Example-Based Slider (E) to express how similar/different the AI-generated notes should be to an example selection, or audition Multiple Alternatives (F) generated by the AI: users select a sample thumbnail to temporarily substitute it into the music score (shown as glowing notes in this figure (G)), then choose to keep it or go back to their original. Users can also use the Infill Mask (H) to crop a section of notes to be infilled again using AI.**

Beyond the infill mask, Cococo distinguishes itself with its *AI steering tools*, which enable novice users to iteratively steer the music co-creation process in real-time. These tools include: 1) *Voice Lanes*, which allow users to progressively define when and where the AI generates music, before any music is created, 2) an *Example-based Slider* for expressing how much the AI-generated music should be similar or different to an existing example of music, 3) *Semantic Sliders* that users can adjust to direct the music towards high-level directions (e.g. happier / sadder, or more conventional / more surprising), and 4) *Multiple Alternatives* for the user to select between a variety of AI-generated options.

**Voice Lanes**

Voice Lanes within Cococo allows a user to specify which voice(s) for which to generate music in a given temporal range. With this capability, users can control the amount of generated content they would like to work with. This was designed to address information overload caused by Coconet's default capabilities to infill all remaining voices and sections at a time. For example, a user can request the AI to add a single accompanying bass line to their melody by highlighting the bass (bottom) voice lane for the duration of the melody, prior to clicking the generate button (see Figure 1b). To support this type of request, we pass a custom generation mask to the Coconet model including only the user-selected voices and time-slices to be generated.

**Semantic Sliders**

We implemented two semantic sliders in Cococo to influence what the generative DNN creates: a conventional vs surprising slider, and a major (happy) vs minor (sad) slider. This was based on formative observations that users wanted to control both musical qualities (e.g., how much the musical phrase or motif should stand out from what is already there) and emotional qualities (e.g., should the notes together produce happy or sad tones).

Users can make the generated notes more predictable given the current context by specifying more "conventional" on the slider, or more unusual by specifying more "surprising." The conventional/surprising slider adjusts the "peaky-ness" of the probability distribution, by mapping to a parameter more formally known as the *temperature T* of the sampling distribution [39]. In formative testing, we found that a log scale interval of $[1/8, 2]$ with a midpoint of $1/2$ yielded a reasonable range of results. In addition, we refined the semantic labels of conventional/surprising based on user feedback to best capture its behavior.

The major vs. minor slider allows users to direct the AI to generate note combinations with a happier (major) quality or a sadder (minor) quality. The limits of this slider include happy and sad face emojis to signal the emotional tones users can expect to control. To generate a passage that follows a more major or minor tone, we define a soft prior (described below) that encourages the sampling distribution to generate the most-likely major triad (for happy) or non-major triad (for sad) at each time-step.

**Audition Multiple Alternatives**

Cococo provides affordances for auditioning multiple alternatives generated by the AI. This capability was designed based

on formative feedback, in which users wanted a way to cycle through several generated suggestions to decide which was the most desirable. We allow the user to select the number of alternatives to be generated and displayed (with a default of three). A thumbnail preview of each alternative is displayed and can be selected for audition within the editor, allowing the user to hear it within the larger musical context. The musical chunk used as a prior to a generation is accessible via the top thumbnail preview (labeled "original") so that users can always compare what the previous version of the piece sounded like, and opt to not use any of the generated alternatives.

**Example-based Slider**
While prototyping the Multiple Alternatives feature, we found that the non-determinism inherent in a deep generative model like Coconet can lead to two opposite, but undesirable extremes: generated samples can be too random and unfocused, or they can be too similar to each other and lack diversity. For example, when the generation area was small relative to surrounding context, generated results would become repetitive: there were a limited set of likely notes for this context according to the model. As a solution, we developed the example-based slider for expressing that the AI-generated music should be more or less like an existing example of music. Before this slider is enabled, the user must select a reference example chunk of notes, either by using the most recent set of notes generated by AI, or manually selecting a reference pattern using the voice lanes or infill mask. Example-based sliders also use soft priors to guide music generation.

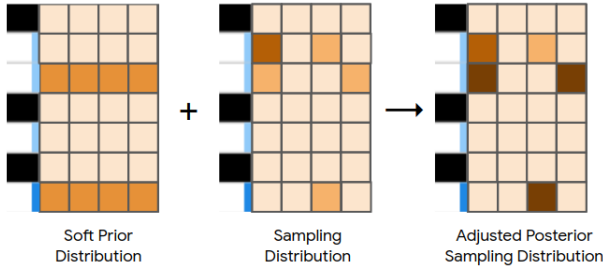**Soft Priors: a Technique for AI-Steering**



Figure 2. Visualization of using soft priors to adjust a model's sampling distribution. The shape of the distributions are simplified to 1 voice, 7 pitches, and 4 timesteps. In CoCoCo, the actual shape is 4 voices, 46 pitches, and 32 timesteps

Many of our AI-steering tools make use of a "soft prior" to modulate the model's generated output. These priors enable users or AI-steering tool designer to add control to existing generative models without needing to retrain them. The model's sampling distribution is a softmax [16] probability distribution over all possible pitches, for each voice and for each time step; high probabilities are assigned to the pitches that are likely given the infill's surrounding musical context. The soft prior approach enables the generation of output that adheres to *both* the surrounding context (encoded in the model's sampling distribution) and additional desired qualities (encoded in a prior distribution). More formally, we use the equation below to alter the distribution used to generate outputs:

$$p_{\text{adjusted}}(x_{v,t}|x_C) \propto p_{\text{coconet}}(x_{v,t}|x_C)\, p_{\text{softprior}}(x_{v,t})$$

where $p_{\text{coconet}}(x_{v,t}|x_C)$ gives the sampling distribution over pitches for voice $v$ at time $t$ from Coconet given musical context $x_C$ ($C$ gives the set of $v,t$ positions constituting the context), $p_{\text{softprior}}(x_{v,t})$ encodes the distribution over pitches specified by the user or AI-steering tool designer (serving as soft priors), and $p_{\text{adjusted}}(x_{v,t}|x_C)$ gives the resulting adjusted posterior sampling distribution over pitches.

The soft priors $p_{\text{softprior}}(x_{v,t})$ are defined so that notes that should be encouraged are given a higher probability, and those discouraged are given a lower, but non-zero probability. This setup allows for two desirable properties. First, since none of the note probabilities are forced to zero, very probable notes in the model's original sampling distribution can still be likely after incorporating the priors. Second, even though the priors are specified for particular voice and time steps, their effects can propagate to other parts of the piece. For example, as Coconet fills in the music, it will try to generate transitions that go smoothly between parts with a soft prior and parts without. These together make it possible for the model's output to adhere to both the original context and the additional user-desired qualities.

The soft priors technique powers Cococo's example-based slider and semantic-sliders. When the user sets the example-based slider to more "similar," we create a soft prior that has higher probabilities for notes in the example. Conversely, for a slider setting of more "different," we create a soft prior that has lower probabilities for notes in the example. The soft prior is then used to alter the sampling distribution according to Equation 1 and Figure 2.

The minor/major slider uses a slightly more complicated approach to define the soft prior distribution. To encourage notes from a major chord, for example, we construct the soft prior by asking what is the most likely major chord at each time-slice within the model's sampling distribution. The log likelihood of a chord is computed by summing the log probability of all the notes that could be part of the chord (e.g., for C major chord, this includes all the Cs, Es, and Gs in all octaves). We repeat this procedure for all possible major chords to determine which chord is the most likely for a time slice. We then repeat this procedure for all time-slices to be generated, in order to create our soft prior for most-likely major chords; this soft prior is used to alter the sampling distribution to create the adjusted posterior sampling distribution as shown in Figure 2.

The features described above were implemented as a React.js web application, backed by an open-source browser-based implementation [35] of the Coconet model. We modified Coconet to include soft priors.

**USER STUDY**
We conducted a user study to evaluate the extent to which AI-steering tools support user needs, and to uncover how they affect the user experience of co-creating with AI. To this end, we compared the experiences of music novices using Cococo to that of a conventional interface that mirrors current interfaces for deep generative models (e.g. the Bach Doodle [25]). The conventional interface is aesthetically similar to Cococo, but does not contain the AI-steering tools. To mirror the most
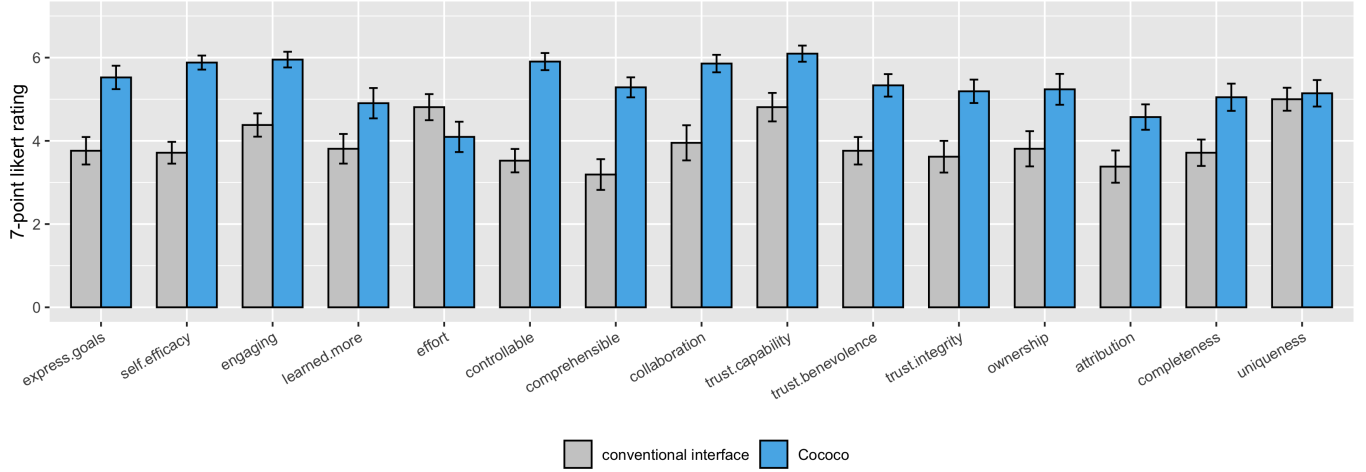
**Figure 3. Results from post-study survey comparing the conventional interface and Cococo, with standard error bars.**

recent deep generative music interfaces, the conventional interface does include the *infill-mask* feature described above, which enables users to crop any region of the music and request that it be filled in by the AI [7, 19]. We ask in this study: **RQ1** How do the AI-steering tools affect user perceptions of the creative process and the creative artifacts made with the AI (e.g., perceptions of ownership, self-efficacy, trust in the AI, quality of the composition, etc.) **RQ2** How do music novices apply the AI-steering tools in their creative process? What patterns of use and strategies arise?

### Measures

To answer the research questions above, we evaluated the following outcome metrics. All items below were rated on a 7-point Likert scale (1=Strongly disagree, 7=Strongly agree, except where noted below).

The following set of metrics sought to measure users' compositional experience. **Creative expression**: Users rated *"I was able to express my creative goals in the composition made using [System X]."* **Self-efficacy**: users answered two items from the Generalized Self-Efficacy scale [36] that were rephrased for music composition. **Effort**: users answered the effort question of the NASA-TLX [20], where 1=very low and 7=very high. **Engaging**: users rated *"Using [System X] felt engaging."* **Learning**: users rated *"After using [System X], I learned more about music composition than I knew previously."* **Quality of the composition**: users rated the *completeness* of the composition and the *uniqueness* of the composition.

In addition, we evaluated users' attitudes towards the AI. **AI interaction issues**: users rated the extent to which the system felt *comprehensible* and *controllable*, two key challenges of human-AI interaction raised in prior work on DNNs [32]. **Trust**: participants rated the system along Mayer's dimensions of trust [31]: capability, benevolence, and integrity. **Ownership**: users rated two questions, one on ownership (*"I felt the composition created was mine."*), and one on attribution (*"The music created using [System X] was 1=totally due to the system's contributions, 7=totally due to my contributions."*).

**Collaboration**: users rated *"I felt like I was collaborating with the system."*

In the study, we called the two systems "System 1" and "System 2" (counterbalanced) to avoid biasing participants, but we refer to them here as "Cococo" and "conventional interface" for clarity.

### Method

21 users participated in the user study. To ensure that participants were novices in composition, we required that participants had played a musical instrument before at some point in their life, but had none or relatively little experience with composition and music theory. Participants were recruited through mailing lists at our institution, and received $40 worth of gift credit for their time.

Each participant first completed an online tutorial of the two interfaces on their own (30 minutes). Then, they composed two pieces, one with Cococo and one with the conventional interface, with the order counterbalanced (15 minutes each). As a prompt, users were provided a set of images from the card game Dixit [38] and were asked to compose music for one that reflected the character and mood of the image. This task is similar to image-based tasks used in prior music studies [24]. Finally, they answered a post-study questionnaire and completed a semi-structured interview (20 minutes). So that we could understand their thought process, users were encouraged to think aloud while composing.

The participants who completed the study included 12 females and 9 males, ages 20 to 52 ($\mu = 31$). Almost all had either very little experience with music theory (12 participants) or a beginner-level understanding of concepts like note reading, major and minor scales and key signatures, intervals, triads and time signatures (8 participants). Participants had diverse prior experiences pursuing music composition, where 6 had never considered composing, 8 had considered composing but never done it, and 7 had tried improvising or creating music informally (e.g., by fiddling around on a piano, making

6

small songs on music software, or creating small music-group arrangements).

**QUANTITATIVE FINDINGS**

Results from the post-study questionnaire are shown in Figure 3. In the post-survey, participants felt they were better **able to express their creative goals** in the composition made using Cococo ($\mu = 5.5$) compared to the conventional interface ($\mu = 3.8$). Results from a paired t-test show a statistically significant difference ($p = 0.0006, t = 4.0$). We found **self-efficacy** was higher when using Cococo than the conventional interface ($\mu = 5.9, \mu = 3.7, p < 0.0001, t = 7.8$), which was obtained by computing an average of our two self-efficacy items (Cronbach's $\alpha = 0.86$). Participants felt Cococo was more **engaging** compared to the conventional interface ($\mu = 6.0, \mu = 4.4, p = 0.0001, t = 5.0$). They thought they **learned more** about music composition after using Cococo compared to after using the conventional interface ($\mu = 4.9, \mu = 3.8, p = 0.0003, t = 4.4$). This positive result suggests that Cococo's steering tools provided more structured ways to engage and learn about music composition, a result we explore further in the qualitative results. No significant difference was found in **effort** between Cococo and the conventional interface ($\mu = 4.1, \mu = 4.8, p = 0.1514, t = 1.5$); participants described the two systems as requiring different kinds of effort: while Cococo required users to think and interact with the controls, the conventional interface's lack of controls made it effortful to express creative goals.

After using the two AI composing tools, participants found Cococo to be more **controllable** than the conventional interface ($\mu = 5.9, \mu = 3.5, p < 0.0001, t = 7.1$) as well as more **comprehensible** ($\mu = 5.3, \mu = 3.2, p < 0.0001, t = 5.6$). Users felt like they were **collaborating** more with the AI when using Cococo compared to the conventional interface ($\mu = 5.9, \mu = 4.0, p = 0.0002, t = 4.5$). Participants expressed higher **trust** in Cococo compared to the conventional interface, along the capability dimension ($\mu = 6.1, \mu = 4.8, p = 0.0008, t = 4.0$), benevolence dimension ($\mu = 5.3, \mu = 3.8, p = 0.0004, t = 4.3$), and integrity dimension ($\mu = 5.2, \mu = 3.6, p = 0.0055, t = 3.1$).

Users felt more **ownership** over the compositions made using Cococo as compared to the conventional interface ($\mu = 5.2, \mu = 3.8, p = 0.0071, t = 3.0$). They felt the music was more **due to their own contributions relative to the AI** when it was co-created with Cococo compared to the conventional interface ($\mu = 4.6, \mu = 3.4, p = 0.0136, t = 2.7$). Participants thought the composition created using Cococo was more **complete** ($\mu = 5.0, \mu = 3.7, p = 0.0116, t = 2.8$). We found no statistically significant difference between participants' perceptions regarding the **uniqueness** of compositions ($\mu = 5.1, \mu = 5.0, p = 0.6507, t = 0.5$).

**QUALITATIVE FINDINGS**

In this section, we describe participants' strategies for co-creating music, how they leveraged the AI-steering tools to work around perceived limitations of the AI, and how the tools helped novices "up-level" their existing skills and knowledge, while still retaining a sense of agency and ownership.



**Figure 4. Common Patterns of using Voice Lanes, visualized using interaction data from 4 archetypal participants (darker-colored segments were performed by users before lighter-colored segments): (A) Voice-by-voice (most common), (B) Temporal Chunks, (C) Combination of Voice-by-Voice and Temporal Chunks, and (D) Ad-hoc Bits**

**Tool-Based Strategies for Composing with AI**

Participants made use of AI-steering tools to support initial brainstorming, to generate alternatives, and to steer the AI when the music did not match their general intent (e.g. mood) or the specific tune in their head. Their overall composition process typically involved breaking the task down into smaller, semantically meaningful pieces, and generating, auditioning and editing content in these smaller pieces until they arrived at a satisfactory result. The AI-steering tools played key roles throughout this process.

*Building Up, Bit-by-Bit*

Many participants used the Voice Lanes to iteratively develop one voice layer at a time, in a "brick-building" fashion (Figure 4a): *"I'm trying to get the bass right, then the tenor right, then soprano and alto right, and build bit-by-bit"* (P2). This use of the Voice Lanes helped reduce the mental workload of handling multiple voices at once: *"As someone who cannot be thinking about all 4 voices at the same time, it's so helpful to generate one at a time"* (P2). Other participants leveraged the temporal aspect of the lanes (Figure 4b), using the AI to generate all four voices in a temporal region then refining the result. Some tried a combination of the voice-wise and temporal approaches, by working voice-wise in the first half of the song, then letting the AI continue a whole temporal chunk in the second half (Figure 4c).

One participant referred to this piece-wise process as creating intermediate "checkpoints," where they stopped and evaluated the song before more content was generated. This strategy allowed participants to *"intervene after [the AI] generated [content]... stop it in the middle... and change it to feel different, before it kept going"* (P14).

In contrast, in the conventional interface, the AI fully auto-completed the music at once. As a result, participants resorted to "sculpting" and refining the AI's fully-generated music by repeatedly using the Infill Mask. Echoing the results in our need-finding study, some participants found the amount of resultant content overwhelming.

*Working With Semantically Meaningful Chunks*

In a similar spirit of composing bit-by-bit, participants actively leveraged the AI-steering tools to divide the music into

semantically meaningful chunks, based on voice or time. For example, many used Voice Lanes to differentiate between the melody and background by using separate voices, or they assigned different musical personas to different voices. For example, one participant gave the tenor voice an *"alternating [pitch] pattern"* to express indecision in the main melody, then gave other voices *"mysterious... dinging sounds"* as a harmonic backdrop (P4).

Participants also divided the music into temporally distinct chunks as a way of illustrating evolution or change. One participant communicated a fight was about to start by requesting more conventional chords in the beginning third of the piece, then used the minor and surprising slider to generate an unresolved feeling in this evolving battle scene in the middle of the piece. In the final section, they used *"prolonged notes [to match] the long stare"* between dueling characters.

### Generating, Auditioning, and Editing
Participants often employed the AI-steering tools to point the AI in a desired initial direction, audition the generated content, or edit and steer the generated output. The Multiple Alternatives functionality naturally lent itself to this "generate and audition" strategy of music composition. Participants could generate a range of possibilities, audition them, and choose the one closest to their goal before continuing.

When generating content, the Semantic Sliders were sometimes used to set an initial trajectory for generated music: *"There's one... idea in my head.... that's the signal that I'm giving to the computer"* (P3). Some felt that this capability helped constrain the large space of possibilities that could be generated: *"Because I was able to give more inputs to [Cococo] about what my goals were, it was able to create some things that gave me a starting point"* (P8). In analysis of logs, 12 of the 21 participants modified the default values of the slider parameters prior to their first generation request to Cococo.

The AI-steering tools were also used to refine what the AI had generated, pushing the content in a direction closer to their intentions: *"It was... not dramatic enough. Moving the slider to more surprising, and more minor added more drama at the end"* (P5). Applying the example-based slider, participants moved the setting to "similar" to push content closer to an example that embodied their musical goals: *"Work your magic on these notes, but keep it similar so they won't move around too much"* (P1). They set the slider to "different" when the initial AI-generated notes were *"not sounding good"* (P15) or when all the generated options needed to be *"totally scrapped"* (P13) because all were of opposite quality to the sound the user desired.

## Tool-based Strategies for Addressing AI Limitations
In this section, we describe ways in which the tools were used to discover and directly address AI limitations.

### Identifying and Debugging Problematic AI Output
By building up the music bit-by-bit, users became familiar with their own composition during the creation process, which enabled them to more quickly identify the *"cause"* of problematic areas later on. For example, one participant indicated that *"[because] I had built [each voice] independently and listened to them individually,"* this helped them *"understand what is coming from where"* (P7). Conversely, if multiple voices were generated simultaneously, participants found it difficult to understand the complex interactions: *"It's harder to disentangle what change caused what... when I make a change, there could be this mixed reaction...it propagates to [multiple] things at once"* (P6). By enabling users to generate bit-by-bit from the ground up, and incrementally evaluate the music along the way, the tools may have enabled novices to better understand and subsequently "debug" their own musical creations.

### Testing and Discovering the Limits of the AI
The tools also enabled participants to discover the limits of the AI. One participant, while using Voice Lanes to generate multiple alternatives for a single-voice harmony, discovered that the AI may be constrained by what's musically possible: *"Maybe the dissonance is happening because of how I had the soprano and bass... which are limiting it... so it's hard to find something that works"* (P15). Here, the Voice Lanes helped this user consider the limits imposed by a specific voice component, enabling them to reflect on the limits of the AI in a more semantically meaningful way. The Multiple Alternatives capability further enabled this participant to systematically infer that this particular setting was unlikely to produce better results through the observation of multiple poor results produced by the AI.

Some participants also set the sliders to their outer limits to test the boundaries of AI output. For example, one participant moved a slider to the "similar" extreme, then incrementally backed it off to understand what to expect at various levels of the slider: *"On the far end of similar, I got four identical generations, and now I'm almost at the middle now, and it's making such subtle adjustments"* (P18). These interactive adjustments allowed the user to quickly explore the limits of what they can expect the AI tools to generate, aiding construction of a mental model of the AI's capabilities. In contrast, when using the conventional interface, participants could not as easily discern whether undesirable model outputs were due to AI limits, or a simple luck of the draw.

### Proxy Controls
Participants drew upon a common set of composition strategies to achieve desired outcomes. For example, higher pitches were used to communicate a light mood, long notes to convey calmness or drawn-out emotions, and a shape of ascending pitches to communicate triumph and escalation.

During the study, participants who could not find an explicit way to express these concepts to the AI re-purposed the AI-steering tools as "proxy controls" to enact these strategies. For example, users sometimes hoped that the surprising vs. conventional slider would be correlated with note density and tempo. A common pattern was to set the slider to "conventional" to generate music that was *"not super fast... not a strong musical intensity"* (P9), and to "surprising" for generating *"shorter notes... to add more interest"* (P15). Participants

also turned to heuristics (such as knowledge that bass lines in music tend to contain lower pitches) to "reverse-engineer" which Voice Lanes to select in an attempt to control pitch range. Multiple steering tools were also sometimes combined to achieve a desired effect, such as using "conventional" in conjunction with the bass Voice Lane to create slow and steady music.

In some cases, even use of the AI-steering tools did not succeed in generating the desired quality. For example, the music produced using the "similar" setting was not always similar along the user-envisioned dimension, and the surprising slider did not systematically map to note density, despite being correlated. Facing these challenges, participants developed a strategy of "leading by example" by populating surrounding context with the type of content they desired from the AI. For instance, one participant manually drew an ascending pattern in the first half of the alto voice, in the hopes that the AI would continue the ascending pattern in the second half.

**Novice Up-Leveling, Agency, and Collaboration**

Beyond assisting with content generation and editing, the AI-steering tools seemed to help participants extend their music composition knowledge and skills.

*Learning and Discovering Musical Structure*

In the Cococo interface, there is no way to request initial music generation by the AI without first selecting Voice Lanes. As a result, the tools implicitly created a more structured workflow, which seemed to be helpful in providing scaffolding for novices: *"With all the controls, I feel more secure..... you have the bars of the [Voice Lanes]... you feel surrounded by this support of the machine"* (P13).

Some participants had difficulty understanding how individual elements of the composition interacted with one another, and re-purposed the AI-steering tools to help improve their understanding. For example, one participant described how a workflow of 1) manually composing a seed voice, 2) using the AI to generate a single accompanying voice from that seed, and 3) modifying the seed and repeating this process helped them *"more directly see how the changes [they] made affect things"* (P6). Another participant was *"curious what [Cococo] will put in for alto...[After the alto is generated] it seems to go with the soprano, but there's some dissonance near the beginning"* (P15). By isolating and revealing the effects of a single voice on another, the tools allowed participants to "micro-evaluate" the music and discover patterns in how components interact.

The tools also helped participants learn how sub-components affect the whole piece. One participant described how they came to understand *"that having that soprano up [at this bar]... gives a total injection of a different emotion,"* which they only realized by using the Voice Lanes to place a single voice within a single bar. Another participant learned that *"a piece can become more vivid by adding both a minor and major chord"* after they applied the major/minor slider to generate two contrasting, side-by-side chunks (P12). Thus, while the conventional AI could do everything on its own, partitioning the AI's capabilities into smaller, semantically meaningful

tools helped people learn music composition strategies that they could re-use in the future.

*Novice Self-Efficacy vis-a-vis the AI*

In post-study interviews, novices described how the tools instilled a sense of competence, self-efficacy, and agency when composing. For example, a participant contrasted the conventional interface, in which the *"machine is doing all the work,"* to Cococo, where they felt *"more useful as a composer"* (P3). The AI-steering tools also seemed to instill a sense of creative agency. By enabling participants to indicate what type of music was generated, the slider controls *"really help to express [myself] in a way [I] wouldn't be able to do in music notes or words"* (P7). Participants also attributed their sense of agency and ownership to the availability of choice, even if it wasn't exercised: *"There are options, but I don't feel like I have to use them... it's not like the [AI] is telling me 'This is the correct thing to do here'... so I felt I definitely had ownership in the music"* (P9). In contrast, participants indicated that they felt less ownership of the music in the conventional interface because they performed a smaller portion of the work, relative to the AI: *"The more I used the AI... the less I personally compose, the less ownership I felt....I was not as creative, I felt like I got lazier with the music...I relied on the AI to solve problems"* (P9).

While there were indications that the tools helped improve feelings of self-efficacy, there were also times when participants questioned their own musical capabilities when they were unable to obtain desirable results. Because the AI generates music given a surrounding "seed" context, users who were dissatisfied with AI output often wondered whether they had provided a low-quality seed, leading to suboptimal AI output: *"All the things it's generating sound sad, so it's probably me because of what I generated"* (P11). In cases such as this, participants seemed unable to disambiguate between AI failures and their own compositional flaws, and placed the blame on themselves.

In other instances, novices were hesitant to interfere with the AI music generation process. For instance, some assumed that the AI's global optimization would create better output than their own local control of sub-units: *"Instead of doing [the voice lanes] one by one, I thought that the AI would know how to combine all these three [voices] in a way that would sound good"* (P1). While editing content, others were worried that making local changes could interfere with the AI's global optimization and possibly *"mess the whole thing up"* (P3). In these cases, an incomplete mental model of how the system functions seemed to discourage experimentation and their sense of self-efficacy.

*Novice Perceptions of AI's Collaborative Role*

The ability to use AI-steering tools also affected how users perceived of the AI as a collaborator. When using Cococo, users conceived of the AI as a collaborator that could not only inspire, but also revise and adjust to requests. For instance, one described it as a nimble team who *"could be adjusted to do what I would like for them to do... I had a creative team [if I needed one] or I had a conventional team [if I needed one]... like a large set of collaborators"* (P19). Others

9

appreciated that Cococo was able to yield control to the end-user, and viewed the AI as more of a highly-proficient helper: *"An art assistant, who is extremely proficient, but has a clear understanding of who is in control of the situation"* (P18).

In contrast, participants perceived the conventional interface differently, calling it a *"brilliant composer"* (P16) they could outsource work to, but who was more difficult to communicate with. When working with the conventional interface, users were optimistic about its ability to surprise them with musical suggestions that they would not have thought of on their own but pessimistic about its *"blackbox"* (P19) persona when communicating and *"take-it-or-leave-it"* (P6) attitude when working together.

These differing views of the co-creation process with the two interfaces led to distinct ideas of where each interface would be most useful. For the conventional interface, participants imagined it to be useful when they feel *"lazy, and need to generate ideas quickly,"* (P2) or when they feel competent to compose most of a piece manually but are open to brilliant, unexpected suggestions. On the other hand, Cococo seemed useful when the user *"has some [creative goals] in mind that [they] want to build upon"* (P13).

## DISCUSSION/IMPLICATIONS

### Onboarding and Increasing AI Transparency

While participants were able to develop productive strategies using AI-steering tools, they were sometimes hesitant to make local edits for fear of adversely affecting the AI's global optimization. These reactions suggest that participants could benefit from a more accurate mental model of the AI. Previous research suggests benefits of educating users about the AI and its capabilities [1], or providing onboarding materials and exercises [3]. For example, an onboarding tutorial could demonstrate contexts in which the AI can easily generate content, and situations where it is unable to function well. For example, the system could automatically detect if the AI is overly constrained and unable produce a wide variety content, and display a warning sign on the tool icon. Or, semantic sliders could divulge certain variables they are correlated with but not systematically mapped to, to set proper expectations when users leverage them as proxies. This could help users better debug the AI when it produces undesirable results. It could also prevent them from incorrectly attributing themselves and their lack of experience in composing as the source of the error, rather than the AI being overly constrained.

### Bridging User Strategies with the AI

Though we created an initial set of AI-steering tools, we were surprised to discover that participants were *already* prepared with their own set of go-to building blocks, including basic concepts such as pitch, note density, and shape, and semantic concepts such as voice-wise separation of foreground vs. background, or temporal separation of tension vs. resolution. When they could not directly enact these strategies, they re-purposed the existing tools to achieve the desired effect. Given this, one could imagine directly supporting these common go-to strategies. Given a wide range of possible semantic levers, and the technical challenges of exposing these dimensions in

DNNs, model creators should at minimum prioritize exposing dimensions that are the most commonly relied upon (pitch, note density, shape, voice and temporal segmentation). Further, our study found evidence that novices may benefit from learning about composition through tool interaction. Future systems could help boost the effectiveness of novice strategies by helping them bridge between their building blocks to high-level creative goals, such as automatically "upgrading" a series of plodding bass line notes to create a foreboding melody.

### Co-Creation with "Super-Human" Classes of AI

It has long been a goal to produce AI that creates as well as or better than humans. The arrival of systems that start to meet this objective enables focused research on how to expose these super-human capabilities for effective co-creation. Our results suggest the value in partitioning these super-human AIs into smaller, semantically meaningful AI-steering tools.

One unexpected side effect was that novices quickly became familiar with their own creations through composing bit-by-bit, which later helped them debug problematic areas. Interacting through semantically meaningful tools also helped them learn more about music composition and effective strategies for achieving particular outcomes (e.g., the effect of a minor key in the composition). Ultimately, AI-steering tools affected participants' sense of artistic ownership and competence as amateur composers, through an improved ability to express creative intent. Though we didn't measure objective quality of output, users also felt their compositions were more complete (see Quantitative Findings) when the tools were available. In sum, *beyond* reducing information overload, AI-steering tools may be fundamental to one's notion of being a creator, while opening the door for users to learn effective strategies for creating in that domain. We expect one would observe similar results for AI-steering tools developed for super-human AIs in other domains.

Our work also uncovers the dual challenges and opportunities of sophisticated DNNs: although such models can be difficult to decompose, they also expose a flexible space for modification. We found the use of "soft priors" to be a relatively lightweight method for nudging the AI's output without retraining the model. This particular technical approach is likely to be applicable to human-AI co-creation tooling in domains where a probability sampling distribution is exposeable from a deep generative model.

### Defining the Human-AI Partnership

Participants' diverse conceptions of the AI's collaborative role raises the question of what it means to co-create with AI, and what constitutes a truly creative partnership. Users perceived of the AI as a responsive collaborator when AI-steering tools were available, whereas when they were absent, participants felt they were merely outsourcing work to a "brilliant composer." Yet, as indicated, some could conceive using the different collaborator personas for different use cases.

Given this, future interfaces might empower users to define the creative objective depending on their current creative mindset, with the human-AI interface adjusting accordingly. For example, when creative goals are fuzzy and flexible, the AI

could encourage ideation by exploring several points in the space automatically. Alternatively, when the user has a clear goal, the AI could adjust the direction of exploration based on more explicit requests. As such, AI-steering tools could be leveraged not only to control the AI's creative direction, but also to explicitly *cede* control when more serendipitous, whimsical encounters are desired.

## CONCLUSION

We found that AI-steering tools not only enabled users to better express musical intent, but also had an important effect on users' creative ownership and self-efficacy vis-a-vis the AI. Future systems should expose mid-level building blocks, divulge the AI's capabilities and limitations, and empower the user to define the partnership balance. Taken together, this work advances the frontier of human-AI co-creation interfaces, leveraging AI to enrich, rather than replace, human creativity.

## REFERENCES

[1] Saleema Amershi, Dan Weld, Mihaela Vorvoreanu, Adam Fourney, Besmira Nushi, Penny Collisson, Jina Suh, Shamsi Iqbal, Paul N Bennett, Kori Inkpen, and others. 2019. Guidelines for human-AI interaction. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 3.

[2] Nicolas Boulanger-Lewandowski, Yoshua Bengio, and Pascal Vincent. 2012. Modeling Temporal Dependencies in High-Dimensional Sequences: Application to Polyphonic Music Generation and Transcription. *International Conference on Machine Learning* (2012).

[3] Carrie J Cai, Samantha Winter, David Steiner, Lauren Wilcox, and Michael Terry. 2019. "Hello AI": Uncovering the Onboarding Needs of Medical Practitioners for Human-AI Collaborative Decision-Making. *Proceedings of the ACM on Human-Computer Interaction* CSCW (2019).

[4] Elizabeth Clark, Anne Spencer Ross, Chenhao Tan, Yangfeng Ji, and Noah A Smith. 2018. Creative writing with a machine in the loop: Case studies on slogans and stories. In *23rd International Conference on Intelligent User Interfaces*. ACM, 329–340.

[5] Nicholas Davis, Chih-PIn Hsiao, Kunwar Yashraj Singh, Lisa Li, and Brian Magerko. 2016. Empirically studying participatory sense-making in abstract drawing with a co-creative cognitive agent. In *Proceedings of the 21st International Conference on Intelligent User Interfaces*. ACM, 196–207.

[6] Monica Dinculescu, Jesse Engel, and Adam and Roberts. 2019. MidiMe: Personalizing MusicVAE. (2019). `https://magenta.tensorflow.org/midi-me`

[7] Monica Dinculescu and Cheng-Zhi Anna Huang. 2019. Coucou: An expanded interface for interactive composition with Coconet, through flexible inpainting. (2019). `https://coconet.glitch.me/`

[8] Chris Donahue, Ian Simon, and Sander Dieleman. 2019. Piano Genie. *IUI* (2019).

[9] Douglas Eck and Juergen Schmidhuber. 2002. Finding temporal structure in music: Blues improvisation with LSTM recurrent networks. In *Proceedings of the 12th IEEE Workshop on Neural Networks for Signal Processing*.

[10] Judith E Fan, Monica Dinculescu, and David Ha. 2019. collabdraw: An Environment for Collaborative Sketching with an Artificial Agent. In *Proceedings of the 2019 on Creativity and Cognition*. ACM, 556–561.

[11] Morwaread M Farbood, Egon Pasztor, and Kevin Jennings. 2004. Hyperscore: a graphical sketchpad for novice composers. *IEEE Computer Graphics and Applications* 24, 1 (2004), 50–54.

[12] Rebecca Anne Fiebrink. 2011. Real-time human interaction with supervised learning algorithms for music composition and performance. *PhD dissertation, Princeton University* (2011).

[13] Satoru Fukayama, Kazuyoshi Yoshii, and Masataka Goto. 2013. Chord-sequence-factory: A chord arrangement system modifying factorized chord sequence probabilities. (2013).

[14] Katy Ilonka Gero and Lydia B Chilton. 2019. Metaphoria: An Algorithmic Companion for Metaphor Creation. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 296.

[15] Jon Gillick, Adam Roberts, Jesse Engel, Douglas Eck, and David Bamman. 2019. Learning to Groove with Inverse Sequence Transformations. *arXiv preprint arXiv:1905.06118* (2019).

[16] Ian Goodfellow, Yoshua Bengio, and Aaron Courville. 2016. *Deep learning*. MIT press.

[17] James Granger, Mateo Aviles, Joshua Kirby, Austin Griffin, Johnny Yoon, Raniero Lara-Garduno, and Tracy Hammond. 2018. Lumanote: A Real-Time Interactive Music Composition Assistant.. In *IUI Workshops*.

[18] Matthew Guzdial, Nicholas Liao, Jonathan Chen, Shao-Yu Chen, Shukan Shah, Vishwa Shah, Joshua Reno, Gillian Smith, and Mark O Riedl. 2019. Friend, collaborator, student, manager: How design of an ai-driven game level editor affects creators. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 624.

[19] Gaëtan Hadjeres, François Pachet, and Frank Nielsen. 2017. DeepBach: a Steerable Model for Bach Chorales Generation. In *International Conference on Machine Learning*. 1362–1371.

[20] Sandra G Hart and Lowell E Staveland. 1988. Development of NASA-TLX (Task Load Index): Results of empirical and theoretical research. In *Advances in psychology*. Vol. 52. Elsevier, 139–183.

[21] Curtis Hawthorne, Andriy Stasyuk, Adam Roberts, Ian Simon, Cheng-Zhi Anna Huang, Sander Dieleman, Erich Elsen, Jesse Engel, and Douglas Eck. 2019. Enabling Factorized Piano Music Modeling and Generation with the MAESTRO Dataset. In *International Conference on Learning Representations*.

[22] Cheng-Zhi Anna Huang, Tim Cooijmans, Adam Roberts, Aaron Courville, and Doug Eck. 2017a. Counterpoint by Convolution. In *Proceedings of the International Conference on Music Information Retrieval*.

[23] Cheng-Zhi Anna Huang, Tim Cooijmnas, Adam Roberts, Aaron Courville, and Douglas Eck. 2017b. Counterpoint by Convolution. *ISMIR* (2017).

[24] Cheng-Zhi Anna Huang, David Duvenaud, and Krzysztof Z Gajos. 2016. Chordripple: Recommending chords to help novice composers go beyond the ordinary. In *Proceedings of the 21st International Conference on Intelligent User Interfaces*. ACM, 241–250.

[25] Cheng-Zhi Anna Huang, Curtis Hawthorne, Adam Roberts, Monica Dinculescu, James Wexler, Leon Hong, and Jacob Howcroft. 2019a. The Bach Doodle: Approachable music composition with machine learning at scale. *ISMIR* (2019).

[26] Cheng-Zhi Anna Huang, Ashish Vaswani, Jakob Uszkoreit, Ian Simon, Curtis Hawthorne, Noam Shazeer, Andrew M Dai, Matthew D Hoffman, Monica Dinculescu, and Douglas Eck. 2019b. Music Transformer. In *International Conference on Learning Representations*.

[27] Mikhail Jacob and Brian Magerko. 2015. Interaction-based Authoring for Scalable Co-creative Agents.. In *ICCC*. 236–243.

[28] Pegah Karimi, Mary Lou Maher, Nicholas Davis, and Kazjon Grace. 2019. Deep Learning in a Computational Model for Conceptual Shifts in a Co-Creative Design System. *arXiv preprint arXiv:1906.10188* (2019).

[29] Janin Koch, Andrés Lucero, Lena Hegemann, and Antti Oulasvirta. 2019. May AI?: Design Ideation with Cooperative Contextual Bandits. In *Proceedings of the 2019 CHI Conference on Human Factors in Computing Systems*. ACM, 633.

[30] Feynman Liang. 2016. BachBot: Automatic composition in the style of Bach chorales. *Masters thesis, University of Cambridge* (2016).

[31] Roger C Mayer, James H Davis, and F David Schoorman. 1995. An integrative model of organizational trust. *Academy of management review* 20, 3 (1995), 709–734.

[32] Changhoon Oh, Jungwoo Song, Jinhan Choi, Seonghyeon Kim, Sungwoo Lee, and Bongwon Suh. 2018. I lead, you help but only with enough details: Understanding user experience of co-creation with artificial intelligence. In *Proceedings of the 2018 CHI Conference on Human Factors in Computing Systems*. ACM, 649.

[33] Christine Payne. 2019. MuseNet. (2019). `https://openai.com/blog/musenet`

[34] Adam Roberts, Jesse Engel, Colin Raffel, Curtis Hawthorne, and Douglas Eck. 2018a. A hierarchical latent vector model for learning long-term structure in music. *arXiv preprint arXiv:1803.05428* (2018).

[35] Adam Roberts, Curtis Hawthorne, and Ian Simon. 2018b. Magenta. js: A javascript api for augmenting creativity with deep learning. (2018).

[36] Ralf Schwarzer, Matthias Jerusalem, and others. 1995. Generalized self-efficacy scale. *Measures in health psychology: A user's portfolio. Causal and control beliefs* 1, 1 (1995), 35–37.

[37] Ian Simon, Dan Morris, and Sumit Basu. 2008. MySong: automatic accompaniment generation for vocal melodies. In *Proceedings of the SIGCHI conference on human factors in computing systems*. ACM, 725–734.

[38] Wikipedia contributors. 2019a. Dixit (card game) — Wikipedia, The Free Encyclopedia. `https://en.wikipedia.org/w/index.php?title=Dixit_(card_game)&oldid=908027531`. (2019). [Online; accessed 19-September-2019].

[39] Wikipedia contributors. 2019b. Softmax function — Wikipedia, The Free Encyclopedia. `https://en.wikipedia.org/w/index.php?title=Softmax_function&oldid=915290592`. (2019). [Online; accessed 19-September-2019].