# Introducing Oddio:

A new audio based social network with a seamless, endless feed

Kyle Barron, Dennis Check, Jackie Hu, Rituparna Roy, Natalie Schade, Michael Silvestre

Carnegie Mellon University

**Introduction**

Oddio is a new social media platform for posting audio clips. Oddio fulfills key user needs that are currently unmet by social media. The sheer amount of content posted to social media is constantly increasing, but users have a desire to reduce the amount of time spent looking at a screen. At the same time─ with the rise of fully wireless earbuds, smart speakers and new audio formats such as podcasts─ audio entertainment is the constant background companion for a rapidly increasing number of people. Oddio's harnesses these rising trends in a new audio based social media platform.

Throughout user testing, users wanted to listen to Oddio as a secondary task, using it during times when they were not allowed to use a screen: driving, cooking, working, cleaning, etc. This pointed us in the direction of a hands-free, audio-based interaction system that lets users perform basic social media tasks: liking, saving, responding to posts, etc.

Over the course of the project, we found that many users felt comfortable accomplishing basic tasks in a purely audio CUI. Meanwhile, we also found that a visual interface was needed for users to feel comfortable performing more nuanced interactions like editing one's social graph, adjusting content preferences and creating or editing audio posts. Oddio represents a novel concept for audio social media that could have real-world viability with further refinement.

**Literature Review**

Before designing a prototype for our project, our team did a literature review that was composed of 12 peer-reviewed articles, as well as referring to usability studies from Nielsen Norman Group. To complement our research, we looked into the greater ecosystem of home agents, apps and platforms creating audio-based social media, and analogous services such as Spotify and Twitch to better understand the market.

As our team was unfamiliar with how voice interactions could be sustainable within a social network, our team decided to look at prior research that focused on the practicalities of audio consumption and audio engagement. Without a visual interface to guide users through audio-based technology, onboarding users transforms as "discoverability" becomes compromised in the audio ecosystem (Furqan, 2017). Discoverability was a key pain point when assessing the usability of voice assistants and home agents. In a Nielsen Norman Group assessment of voice assistant usability, they identified two major problems: 1) you have to remember the name of the skills the agent can carry out, and 2) you have to remember the "magic words" that invoke these tasks (Budiu).

Research in VUIs has also explored the capacity for consuming audio content over a long period of time. In an experiment by David McGookin and Stephen Brewster called PULSE, which was an app that converted tweets into whispered conversations, participants found that mundane content was more interesting (Chamberlain, 2017). McGookin and Brewster also uncovered through this app that participants found a quick succession of notifications to be annoying. With this in mind, it became apparent that an audio-based application would require trade-offs between ambient delivery of content, while also allowing for consumers to find compelling content they could react to as well.

An additional hypothesis our team had was built around whether or not voice-centered social media would lead to self-censoring, or even elicit some form of empathy within users in its network. As regret is a common experience when posting content on social media platforms, and a frequent cause of regrettable decisions on Facebook comes from emotional "hot states" or impulsive posts made without reflection (Wang, 2011). Work in this space has been rich particularly within online gaming, as networks such as Xbox Live have been in existence since the early 2000s. Experiences around voice interactions on these platforms seems decidedly mixed, though a large part of this is due to voice communication in games compromising your virtual identity (Carter, 2015).

Beyond looking at prior research within the field, we completed a competitive analysis of the audio-based application landscape. Two major platforms looking to create similar experiences as our app were Riffr and Wavve. Riffr creates a platform that allows people to share 3 minute snippets that you can post to your timeline. However, Riffr relies too closely on the conventional visual model of consuming posts. The timeline is difficult for people to have control over finding posts they want to listen to, and searching is made more challenging by confining audio content to a visual model. Wavve, similarly, is about translating audio experiences into having complementary visual components. Though, as later detailed, we did find a need for visual aids with audio posts from our speed dating interviews, we wanted to focus our research on creating sustainable audio-exclusive interactions that could allow people to engage with friends and content creators in environments that were more conducive to audio, particularly in situations where people might be multitasking or experiencing heavy exhaustion from staring at screens.

## Storyboard & Speed Dating: The Method

In order to understand the needs and desires of users within our projected target user group, team SoundClout went through two rounds of storyboard testing with 10 different participants. Our research goal for the testing was the exploration of how users preferred to interact with audio posts. The interactions that we looked into included comments, editing of audio (such as voice filters), controls including voice, gesture and visual controls, and more.

As mentioned before, we utilized two rounds of storyboarding to test and verify some of our initial hypotheses on how people wanted to consume audio based social media content. Each storyboard was built around a "how might we?" statement. We speed dated the storyboards during each round of testing; the facilitator showed the individual participants each storyboard one at a time. As the participants viewed each storyboard, they were asked to speak about their initial reactions to the story, about whether

they themselves would be interested and comfortable doing what was portrayed in the storyboard. For the first round, we showed six participants thirteen different storyboards. For the second round, we showed four participants five different storyboards.

## Storyboards & Speed Dating: The Result

For the first round of storyboarding we received a variety of positive and negative results. There were extremely positive responses to the idea of funny audio filters. People were also interested in the idea of listening to audio threads as a way to interact with their friends on a social media platform. This idea was predicated on the idea that they would be listening to their friends, not to a transcribed audio read by a CUI (such as alexa or google home). People seemed receptive of, and comfortable with, the idea that some sort of visual would be involved in the platform. Without a visual, people worried that they would make mistakes, and participants expressed a large fear that they would accidentally upload something to the social media platform. The final idea that was positively received by the participants surrounded the ability to share clips of audio/podcast content with friends. In addition to the positive responses to various ideas, we found three areas to which people had strong negative reactions. People did not like the idea of using clapping or gestures to control an audio based social media site. To many, these interactions presented far too large an opportunity for error. There were also concerns about decreased functionality without a visual interface.

The result for our second round of storyboarding also had a variety of positive and negative results. We received positive responses from all of the participants to the idea of funny audio based sound filters. This is an idea that also received positive feedback in the first round. In the second round, we presented ideas that were received well the first time, as well as new ideas. One of the new ideas that was received very positively, is the idea of having integrated visual aids that convert the audio to text. Participants felt that this would increase accessibility and usability during situations where listening to

audio is not feasible. The idea that received the largest amount of negative feedback in this round was about environmental sound needs, the participants largely felt that need was already covered by existing options.

## Survey: The Method

As audio social interaction is not very popularly applied in the market yet, before conducting more detail-oriented qualitative research, our team decided to use a survey to reach out to a larger audience group and to validate some of our broad assumptions. The primary goal of this survey is to understand the scenarios and preferences of current usage of audio consumption, and find if there are any correlations with demographics such as gender, age, and occupation. The survey consists of sixteen questions with a mix of quantitative questions(how many people choose X) and qualitative questions(open-ended). The major questions are designed to explore: 1) the existing habits around current audio listening; 2) possible interactions with audio posts; 3) types of content people would like to listen to.

## Survey: The Result

We received feedback from thirty-seven people in total. Half of them are between 18-24 years old, belonging to the user group we initially planned to target at. The results didn't show strong correlation between audio listening preferences and demographic. We have three major takeaways from the survey results:

Firstly, almost everyone who filled in the survey (97.3%) indicates that they usually multitask when consuming audio contents. It suggests that currently people rarely listen to audio contents just to immerse themselves into the world of sound, but use it as a secondary task to complement something else

they are working on. As for what people are usually doing when listening to audio contents, the top three responses are in commute (86.5%), doing work (75.7%), and cooking (64.9%).

Secondly, people avoid opening audio contents in a public setting when they are not headphone-ready. This comes from the question of *"When you encounter social media content with sound(Instagram video, etc), do you turn on your volume and listen to it?"* We got half and half mixed responses, and when we asked them why or why not in the followed up question, the most frequent answers are associated with public settings and headphones - they don't want to disturb others and don't always have their headphones or earphones on.

Lastly, people listen to music (91.9%), podcasts (70.3%), and meditation narration/music (32.4%) the most. In the meantime, people indicate interests in listening to stand-up comedy, talk shows, and event commentary.

Half of the participants claim that currently, they spend 1-3 hours listening to audio media. Interestingly, half of the people say they usually don't create any contents on social media. The survey shows us people don't have a strong preference on using audio social media hands-free or not yet, and our assumption is that the preference has to do with the infrequent usage and distrust in conversational agents in general.

**Directed Storytelling - The Method**

After getting a broad sense for the types of preferred interactions for audio through the storyboard activity and for people's existing habits and current audio listening practices through the survey, we wanted to understand in depth why people preferred those interactions and practices. We conducted directed storytelling with the research goal to understand why people listen to audio and what are their needs for listening to different types of audio content. We interviewed 6 people from varying backgrounds and ages 22-32 years, and asked them to describe a time when they last heard an audio

content or clip, other than music. The reason to go away from music was to dive deeper into specific use cases of other types of audio content, as we didn't want many music-related components on our platform. The interviews lasted about 20-30 mins and we asked several follow-up questions to understand the whys behind their actions and few questions on audio creation and audio sharing.

**Directed Storytelling - Results**

The results from our directed storytelling included several findings that helped us uncover different types of needs behind listening to audio content. Most people preferred to listen to audio content because they can multitask with it. Most multitasking activities included house chores, exercising or doing manual/clerical work, and being able to do things simultaneously made people feel more productive. While the majority of these people were mostly listening to podcasts, the need behind was also to feel productive with their time, seek new knowledge and discover new content. People wanted to learn about things that interested them the most and especially to discover new types of interesting content aligned with their interests or with specific types of genres, which podcasts offer a lot. When they were asked if they would like to create audio content, most of them positively replied, showing enthusiasm to be able to share their own knowledge and experiences, however, creativity and time constraints were often cited as reasons to not be able to create podcast-like long content. Lastly, people also expressed preferences for sharing audio content with others. Especially when it comes to podcasts, people would share with those who were genuinely interested in similar topics or the content was relevant to them in some ways.

**Discussion: How we informed our pretotype?**

We wanted to understand how people reacted to threads of their friend's audio. This was something that people expressed interest in, but we wanted to understand how an actual thread would

affect people. People were interested in the idea, but without any current example of the idea, people had no experience with which to consider this idea. As such, we started collecting audio clips that were donated by the people within the 2020 MHCI cohort. After putting the thread of audio together using these donated

**Pretotype - The Method**

Based on the results of our previous methods, the team developed a "pretotype" to gauge how users would actually feel about specific elements of our idea. The pretotype had 3 research goals which it aimed to generate answers for. They were the following:

1. Understand how people respond to byte size audio content

2. How would users like to listen to the audio

3. How would users like to interact with the audio

The pretotype itself was a string of audio clips generated from members of the MHCI program at CMU. Each team member asked different members of the program to record a quick ~15-30 second audio clip on their phone of them answering the prompt "What's on your mind?". We then took around 12 audio clips we generated and strung them together in a random order. We then asked people both in and outside of our program to listen to this string of audio clips. We ended up testing this pretotype on 8 different people and asked them basic questions after their experience like "What did you think of these clips? Would you listen to something like this? How would you like to interact with the audio?". We also made sure to contextualize the pretotype by telling participants that what they were about to listen to was "an audio-only social media app, in which their friends were talking about their day".

**Pretype - The Results**

We found several interesting results from testing the pretotype with our participants. Some of the feedback was usability related, and others were conceptually related to the capabilities of the app. **Usability**: Participants indicated that it was difficult to skim audio, they would want to adjust playback speed, and that they would want to be able to alter their voice if they had the chance to make a post. The ability to use this app hands free meant a great deal to participants. The idea of being able to do something like the laundry or cleaning the dishes while listening to social media was very powerful for some of our participants.

**Concepts and Features**: Participants mentioned that they would want to communicate or interact with their friends in some way with the app, rather than just listening to individual clips endlessly; A back-and-forth was necessary to them. We also found that some participants wanted an overview or introduction and visual affordance for audio clips before listening to them. Participants wanted to be able to easily discover new content, as well as use tags or topics which they found relevant to sort posts by. These relevant posts were most likely to be stories by people they know on a personal level over generic posts like "Comedy" or "News". Additional "standard" social media features, such as following, commenting, or liking, were mentioned as wants by our participants.

Taking into account these findings, we discuss their implications for the future of our app's design and functionality.

**Discussion: How we informed our prototype?**

As our pretotype was still quite abstract for people to imagine what the app we were designing would be like, we further consolidated two final variations of prototypes. Our goal was to find out the comfortability level people had with voice commands and what was the minimum visual design we should incorporate. Visually, we added sliders because our research showed that people like to adjust their

feed; for audio experience, we listed a series of voice commands so that people can refer to while engaging with the platform. We also added comment features as we saw people like to interact with audio posts and some more social media based features. Lastly, we added categories for people to understand and get an overview of the type of content they will be listening or posting to.

## Parallel Prototype (Usability Testing) - The Methods

Our prototype aimed to be the best representation of our idea for an audio-only social media platform. We used parallel prototyping to create two different versions of the idea to test. These two different ideas were: 1.  The "audio-only" facing side of the idea, in which participants would have to navigate our app using only voice commands, and 2. The "Audio with Supplementary Visual Stimulus" idea, in which participants were still asked to use voice commands, but could also navigate the app through regular means of tapping or swiping visual elements on a screen.

The way in which we tested our two prototypes was by using the "Wizard of Oz" method of design research. We conducted 8 sessions, choosing a between-groups methodology where 4 participants were asked to complete tasks for the "Audio-Only" interface, and the other 4 were asked to complete the same tasks using the "Audio with Supplementary Visual Stimulus" interface. We then compared feedback across the two different interfaces.

For both of the interfaces, we had 1 of our team members act as a facilitator, and another team member as the "wozer" who would control the prototype. As the facilitator asked the participant to complete tasks, and the participant interacted with the prototype, the "wozer" would change screens, give audio feedback (like saying "now recording, beep beep beep") to simulate what the actual app would do. Tasks included: 1. Search for comedy, save that comedy clip to your favorites, and leave a comment on the audio clip. 2. Listen to your friends audio clip, view/listen to all the comments on that clip, and leave a

comment on this audio clip. 3. Create your own audio clip for your friends to listen to telling them how you are today. 4. Adjust your newsfeed to mostly contain news-related content.

Our parallel prototype and Wizard of Oz methodology aimed to accomplish the following goals:

1. Understand in depth specific needs for interactions

2. Analyze what are some natural and easy to use interactions on this platform

3. Understand if they would use key features of the platform such as listening to posts, commenting on posts, creating posts, adjusting feed

**Parallel Prototype (Usability Testing) - The Results**

The results from our prototype and Wizard of Oz sessions included several findings to help us improve the design and functionality of our app.

Across the four tasks, we found that most users are not used to using voice commands, especially in public spaces; audio-to-visual transcription would be useful. An audio feedback would be less uncanny for the user to talk to. The app feels like Apple Music for podcasts. Using a pin to save was confusing; replace it with a bookmark sign. There was confusion between heart and the save button. Users would probably want to read the comments visually rather than just audio only. Users would like a real time comment feature to leave comments at various timestamps. Many users experienced frustration with having to use audio-only, and that they would prefer to have a screen to interact with. For more detail of additional results, please check the Appendix section at the end of this paper.

Overall these findings greatly helped us determine what our final iteration of the prototype would look like, as well as what key components and functionality it would have. Next we will discuss these design implications and ideas based on our results.

**Discussion - How we informed our final product iteration?**

Findings from our parallel prototype testing were deeply embedded in user preferences and helped us iterate our design to support user journey through the platform and enhance the experience of using the platform. First of all, we decided to include a more consistent tutorial for on-boarding so that users can take their own time to get used to voice commands that they can use to listen to the platform. This is achieved by having a screen dedicated to a list of voice commands to be used which is easily accessible in the app all times. Next, we decided to display transcribed texts for simultaneous reading when listening to audio clips, with an option to listen to raw audio as well when they are in a hands-free mode. We also included three advanced features such as adding real-time audio comments on specific parts of the audio clips, as it allows for targeted experience, adding carefully crafted sound emojis for enhanced experience of interaction, and adding the ability to share snippet of audio so users can share relevant pieces of clips and have a better targeted experience. We also added a slider to allow for control of exploration of content. And lastly, we adjusted some UI and VUI elements that will provide better feedback and feedforward for using the platform. The images of our latest product iteration can be seen in the appendix.

**Next Steps**

As next steps, we would like to test these new features with more people in order to better design for usability of these features which are much required by users, as seen in the previous prototyping activities. Further testing will also help us understand what balance of features are required to make this product adoptable.

**Annotated Bibliographies**

Case, A. (2016). *Calm technology: principles and patterns for non-intrusive design*. Sebastopol,

      CA: OReilly Media.

In this book, Amber Case argues that our attention is overdrawn by technology and we should improve

the relationship between human and digital devices by dissolving technology to our ambient

surroundings. Particularly, Case identifies 8 rules to implement calm technology in chapter 2 "*Principle*

*of Calm Technology*". The principles include "technology should require the smallest possible amount of

attention(I)" and "technology should make use of the periphery(III)". We find these two principles are

especially related to our project, as we intend to design a platform that can be used as a secondary side

task people can do while focusing on something else, or to simply listen when having screen exhaustion

with peripheral attention. Case introduces an attention model that delegated the level of attention as

primary, secondary, and tertiary attention. He argues that both visual and audio only interfaces require the

majority of our limited attention, and calmer technology should present information in parallel. How to

create an audio platform that gives unobtrusive yet unambiguous signals? We want to follow the

principles in this book and find out more.


Hindus, D., Ackerman, M.S., Mainwaring, S.D., & Starr, B. (1996). Thunderwire: a field study of an

      audio-only media space. *CSCW '96.*

In this novel field study, Hindus et al. looked into how a workgroup was able to function using an

audio-only media space known as Thunderwire. While the field study was done way back in 1996, our

team still felt like some of the findings and core idea of the project was extremely relevant to our idea.

Investigations of an audio-only media site concluded that audio-only was suitable for maintaining a media

site, however many adaptations were needed by the users in order to account for some of the nuances and

difficulty of having to rely solely on Audio to interact with each other. This holds very true to some of our findings in which users had a difficult time having to adapt to the audio-only controls, but we believe with experience the users would have a much better time on our app.

Kang, R., Brown, S., & Kiesler, S. (2013). Why do people seek anonymity on the Internet? Informing policy and design. *Conference on Human Factors in Computing Systems - Proceedings.* 2657-2666. 10.1145/2470654.2481368.

In this paper, Kang interviewed individuals all over the world who seek out anonymity online. The purpose was to gain an understanding for why people want to remain anonymous when they browse, as well as some of the behavior which people have while anonymous. We believed this information could help us in our research because of the social media nature of our app. What Kang found was that many past-life experiences of individuals is what pushes them to try and remain anonymous online. An example of this was an interview with a parent, in which they said "When you work with kids, a lot of people feel like you don't have a right to a personal life. You have to be a role model at all times, even when you're not at work." This experience of being a parent made the participant feel as if they needed anonymity in order to live out a personal life online. Our team needed to keep in mind how anonymity might impact the audio-only element of our app, and that voice filtering may be a way to help keep people anonymous. We also realize that anonymity can lead to increased cases of trolling, but solving this issue went beyond the scope of the project.

Myers, C., Furquan, A., Nebolsky, J., Caro, K., Zhu, J. Patterns for How Users Overcome Obstacles in Voice User Interfaces. CHI 2018. Montreal, QC, Canada.

This paper explores the workarounds users face when VUI's NLP systems fail them. Though this paper chose participants that had higher skill levels with technology than average, it's findings found that NLP

errors weren't nearly as big of a threat to UX in VUIs than experience troubleshooting unfamiliar intents or failed feedback (6). With NLP errors, participants were able to troubleshoot what the system identified incorrectly, and hyperarticulate those errors. However with unfamiliar intents, the system was unable to support what the participant needed and it quickly begat more emotions of frustration and confusion (6). This conclusion was pretty informative to our work, because it highlighted that we needed to be explicit in the actions users could carry out while interacting with people in their network.

Vazquez-Alvarez, Y., Brewster S. (2011). Eyes-Free Multitasking: The Effect of Cognitive Load
      on Mobile Spatial Audio Interfaces. CHI 2011. Vancouver, BC, Canada.

This paper explores the cognitive load that is created when using mobile spatial audio interfaces. Vazquez-Alvarez and Brewster identify two areas of difficulty for users when multitasking with audio. Area one is interference amongst multiple audio streams, and area two is the constraints of cognitive load. The results of this research found that people preferred spatial audio techniques (such as ambient music where the listener is not in control of the start and stop) when the cognitive load was low (such as listening to classical music). However, when the cognitive load was high (such as listening to news podcasts), people preferred an interruptible single audio stream, meaning they could turn the audio off and on whenever they wanted.

Lee, Y.-H., & Hsieh, G. (2013). Does slacktivism hurt activism?: The effects of moral balancing and
      consistency in online activism. *Conference on Human Factors in Computing Systems -*
      *Proceedings.* 811–820. 10.1145/2470654.2470770.

This paper explores opportunities of activism online and the findings seemed useful in understanding online user engagement. Understanding what causes users to consistently be motivated to engage was useful and we took inspiration from that finding to design features for creating content on our platform.

Although the paper's focus is very different from our project focus, the methods they used inspired us to test our prototypes parallely as well.
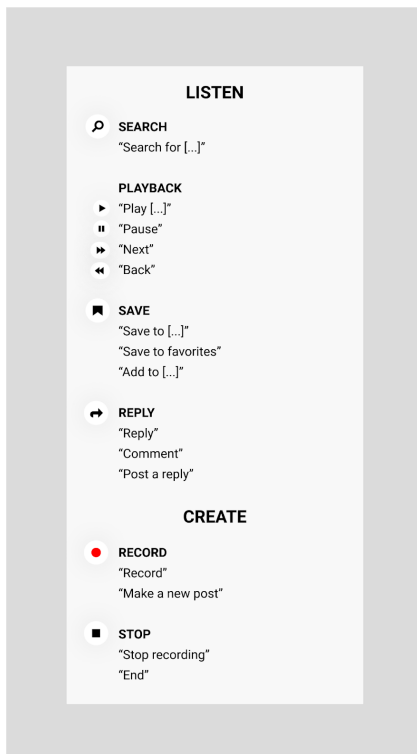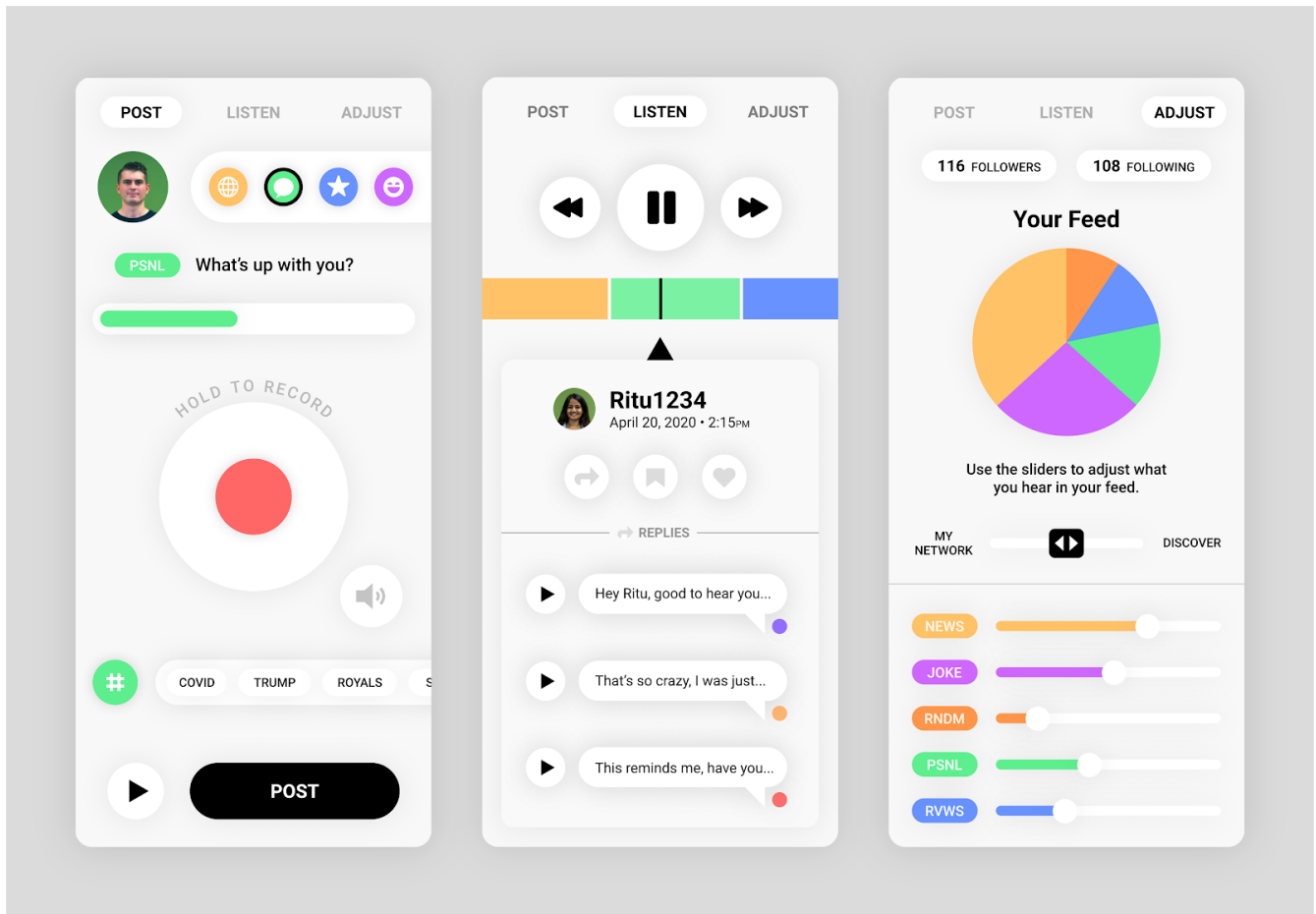
Wang, Y., Komanduri, S., Leon, P., Norcie, G., Acquisti, A., & Cranor, L. (2011). "I regretted the minute I pressed share": A Qualitative Study of Regrets on Facebook. *Proceedings of the Seventh Symposium on Usable Privacy and Security.* 10.1145/2078827.2078841.

This paper deeply explored user responses on social media that leads to guilt and embarrassment, however some of the findings are also particularly useful for us to understand why people take those actions and how we can use that understanding and the downside of that action as an opportunity to design for safer interactions on the platform.
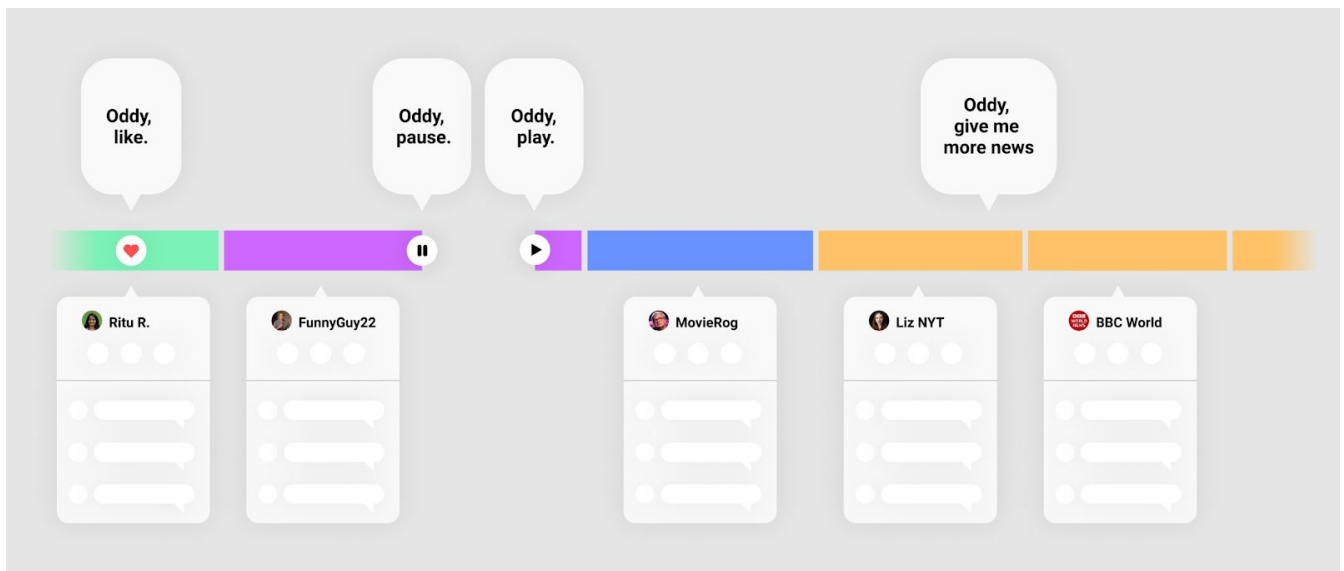
# References

Budiu, R., & Laubheimer, P. (2018, July 22). *Intelligent Assistants Have Poor Usability: A User Study of Alexa, Google Assistant, and Siri*. Retrieved from https://www.nngroup.com/articles/intelligent-assistant-usability/

Carter, M., Gibbs, M., & Wadley, G. (2015). Voice in Virtual Worlds: The Design, Use, and Influence of Voice Chat in Online Play. *Human-Computer Interaction, 30(3-4),* 336–365.

Chamberlain, A., Bødker, M., Hazzard, A., Mcgookin, D., De Roure, D., WIllcox, P., & Papangelis, K. (2017). Audio Technology and Mobile Human Computer Interaction: From Space and Place, to Social Media, Music, Composition and Creation. *International Journal of Mobile Human Computer Interaction (IJMHCI).* 9. 10.4018/IJMHCI.2017100103.

Furqan, Anushay. (2017). Learnability through Adaptive Discovery Tools in Voice User Interfaces. (Publication No. 10601024). [Master's Dissertation, Drexel University]. ProQuest Dissertations Publishing.

# Appendix





**Above: The three main views of Oddio: Post, Listen and Adjust**

**Left: Available CUI commands, shown to users at start of testing**

Track your progress through the stream, or swipe forward and back to navigate quickly

Reply, save and like with a graphical or conversational interface

POST

LISTEN

ADJUST

Ritu1234
April 20, 2020 • 2:15PM

REPLIES

Hey Ritu, good to hear you...

That's so crazy I was just...

This reminds me, have you...

Use buttons to navigate your stream, or voice commands as Oddio runs in the background

See a transcribed preview to help you keep track of comments

Hit play on a comment to pause the main stream and listen to replies

Oddy, like.

Oddy, pause.

Oddy, play.

Oddy, give me more news

Ritu R.

FunnyGuy22

MovieRog

Liz NYT

BBC World

**A conceptual model of the Oddio feed, with CUI-based interaction points at top, and GUI view at bottom**

Select a category for
your post so it finds
the right listeners

Build your clip here—see
how much time is remaining
and add multiple parts

Use hashtags to
help users find
related posts

POST

LISTEN

ADJUST

PSNL

What's up with you?

HOLD TO RECORD

COVID

TRUMP

ROYALS

POST

Add sounds you've
saved or featured sound
effects to your post

When you're done, post
your clip so others can
hear it in their Oddio feeds

Preview the post
you've created
before sharing it

See a quick breakdown
of your feed that
changes as you adjust

Adjust the balance between
familiar and new content with
a bidirectional slider

POST

LISTEN

ADJUST

116 FOLLOWERS

108 FOLLOWING

Your Feed

MY
NETWORK

Use the sliders to adjust what
you hear in your feed!

DISCOVER

NEWS

JOKE

RNDM

PSNL

RVWS

Access and edit
your social graph

Use sliders to adjust the
proportion of each
content type you hear

**Additional Results from Prototype Testing:**

In task 2, we found:

Using the app was initially confusing; It was like whatsapp but prettier. Lack of color makes it hard to navigate. Users wanted to unfold the comments as they listened to the audio piece. They thought it would be easier and faster when they are on screen, unless hands-free. Converting the audio to text, but having the CUI play the raw audio would be cool. Users had the problem of if you can't speak right now, then you can't leave comments. Users would want to have a feature to upload mp3 audio as well. Maybe they can share snippets of audio. Also, sound emoji effects (more like reactions) - need to be very carefully crafted but could enhance the experience.

In task 3, we found:

People were confused about news/personal/reviews categories. Some users would want the default to be personal rather than news. People generally got the concept of hold to record, but there was slight confusion on what they expected when it was just audio. Users didn't notice the top navigation bar, and we should include feedforward for "hold to record". Users wanted to be able to do basic edits + audio auto-tune. Users hated hearing their own voice played back to them after making a recording, so having an ideal way of allowing edits before posting would be necessary and difficult to fix. In order to learn the app, users would want a tutorial of a natural flow which would navigate the visual components, but also on-board them for the audio-only components they could use instead of the visual components. Many users wanted a visual interface to make edits and found the audio-only app difficult to use when making a post.

In task 4, we found:

Being able to curate the content was the most appreciated aspect of our idea. Users thought it was new and innovative. The sliding seemed very intuitive when dictating how much of certain categories of

content would be given on the feed. Users thought that this would solve the problem of their newsfeed

being "ruined" by liking something on accident and then having their feed overrun by that type of content.

A thing to consider is the name "tweak" to something else, as well as including either an "explore"

section in conjunction with the "feed" so that users would not be intimidated to mess with their content.