# Dynamics of Viewer's Attention: Making Content More Engaging in a Vertical Frame

Afzaal Yousaf Baig [1]   Muhammad Riaz [2]   Muhammad Ali Baig [3]   Mirza Zaib Hasan [4]

**ABSTRACT:** The proliferation of mobile-first social media platforms has established vertical video (9:16 aspect ratio) as a dominant format for content consumption. However, established principles of visual composition, largely derived from landscape cinematography, may not directly apply to this narrow, portrait-oriented frame. This study investigates the impact of subject positioning on viewer attention within vertical videos. Using mobile camera-based eye-tracking technology provided by RealEye.io, we conducted a within-subjects experiment with 40 participants. Each participant viewed two purpose-created, 20-second stimulus videos. The first video (Dynamic AOI) featured the primary Area of Interest (AOI) in multiple, varied locations throughout the frame. The second video (Central AOI) consistently positioned the AOI in the central third of the frame. Key eye-tracking metrics were analyzed, including fixation duration, saccade count, and a derived attention efficiency metric, Coefficient K (ratio of fixations to saccades). The results revealed a statistically significant difference between the two conditions. The Central AOI video elicited significantly longer average fixation durations and fewer saccades compared to the Dynamic AOI video. Consequently, Coefficient K was substantially higher for the centrally framed content. These findings suggest that central framing in vertical video reduces cognitive load, minimizes visual search behavior, and facilitates more sustained and focused attention on the primary subject. The study concludes that for content creators, marketers, and platform designers seeking to maximize viewer engagement and message retention in the vertical video ecosystem, a deliberate strategy of centering the primary subject is demonstrably more effective.

**KEYWORDS:** Vertical Video, Eye-Tracking, Visual Attention, Cognitive Load, Fixations, Saccades, Center Framing

[1] PhD Scholar, Faculty of Social Sciences & Humanities, Riphah International University, Islamabad, Pakistan.
Email: afzaal.yousaf@riphah.edu.pk

[2] Assistant Professor, Faculty of Social Sciences & Humanities, Riphah International University, Islamabad, Pakistan.
Email: muhammad.riaz@riphah.edu.pk

[3] Lecturer, Sahara College, Narowal, Punjab, Pakistan.
Email: baigmuhammadali1@gmail.com

[4] Lecturer, Faculty of Social Sciences & Humanities, Riphah International University, Islamabad, Pakistan.
Email: zaib.hassan@riphah.edu.pk

Corresponding Author: Afzaal Yousaf Baig
✉ afzaal.yousaf@riphah.edu.pk

## Introduction

The digital media landscape has undergone a seismic shift in the last decade, characterized by the meteoric rise of mobile devices as the primary screen for information and entertainment consumption. This transition has given birth to a new native language of visual communication: vertical video (Mohammad, 2022). Platforms such as TikTok, Instagram Reels, and YouTube Shorts have amassed billions of users, all engaging with content specifically designed for the 9:16 aspect ratio of a smartphone held in portrait orientation (Cipresso et al., 2018). In 2022, vertical video formats accounted for most of the time spent on these leading social platforms,

making an understanding of their cognitive and perceptual impact not just an academic curiosity, but a critical necessity for effective communication in the 21st century (Álvarez-Álvarez & del Puerto Carrizosa, 2022).

Despite its ubiquity, the "grammar" of effective vertical video is still being written. Traditional filmmaking and photography have long relied on established compositional rules, such as the Rule of Thirds, leading lines, and frame balancing, to guide the viewer's eye and create aesthetically pleasing images (Raursø et al., 2021). These principles were developed and optimized for landscape (horizontal) formats, which mirror the binocular field of human vision and the proscenium arch of the theater. The vertical frame, however, presents a fundamentally different canvas, one that is narrow, tall, and often consumed in rapid succession in a visually cluttered feed. This constrained format may impose unique cognitive demands on the viewer, potentially altering how attention is allocated and sustained (Ohta et al., 2023).

The central challenge for content creators and marketers is to capture and hold attention within the first few seconds of a video. In a scrolling-based interface, viewer attention is scarce and fleeting (Mazzucchi, 2017). A failure to immediately engage the viewer results in a "swipe away," rendering the content unseen and its message lost. Therefore, understanding the elemental building blocks of attention in this format is paramount. One of the most fundamental choices a creator makes is where to place the subject, person, product, or point of action within the frame. Does the Rule of Thirds still hold primacy, or does the narrowness of the vertical format favor a more direct, centered approach?

This study addresses this question directly by investigating how the spatial positioning of the primary Area of Interest (AOI) affects visual attention patterns. We hypothesize that the cognitive processing of vertical video differs from that of traditional landscape media. Specifically, the reduced horizontal space may render off-center positioning less effective, creating a "search task" that increases cognitive load and fragments attention. On the contrary, placing the subject consistently in the center may align more naturally with the viewing habits of mobile users, reducing the need for extensive eye movements and allowing for deeper cognitive engagement with the subject matter (Wolf et al., 2023).

To test this hypothesis, we employed Mobile camera-based eye-tracking technology, a powerful tool for objectively measuring subconscious visual behavior (Kiefer et al., 2017). By analyzing viewers' eye movements, specifically their fixations (periods of sustained gaze) and saccades (rapid movements between fixations), we can quantify the allocation of visual attention and infer the cognitive effort required to process the visual information (Goetz & Neider, 2024).

This research compares two distinct video conditions: one where the AOI is dynamically positioned across different parts of the frame, and another where the AOI remains consistently in the center. Our primary research question is whether the spatial positioning of the primary Area of Interest (AOI) in vertical video significantly affects viewer attention, as measured by fixation duration, saccade frequency, and cognitive efficiency.

The findings of this paper aim to provide empirical, data-driven guidelines for content creators, digital marketers, user experience (UX) designers, and communication scholars. By explaining the relationship between subject positioning and attention, we can move from circumstantial best practices to scientifically validated principles for creating more engaging, effective, and impactful vertical video content.

## Literature Review

The transition from landscape to vertical video is a direct consequence of the smartphone's dominance. As of 2023, mobile devices account for nearly 60% of all web traffic worldwide, and users hold their phones in portrait orientation approximately 94% of the time (StatCounter, 2023; Rathod & Agal, 2023). Platforms that embraced this behavior have thrived. TikTok, launched internationally in 2017, rapidly became a cultural and economic powerhouse, predicated entirely on a full-screen, sound-on, vertical video experience. Its success prompted established players like Meta (with Instagram Reels) and Google (with YouTube Shorts) to pivot their strategies aggressively toward the same format (Stokel-Walker, 2021).

This format is not merely a cropped version of landscape video; it constitutes a distinct medium with its own affordances and constraints. The 9:16 aspect ratio emphasizes height over width, making it well-suited for single-subject portraits, direct-to-camera addresses, and showcasing linear, vertical motion (e.g., a person dancing, a waterfall). However, it is less suited for capturing expansive scenes, group interactions, or horizontal action, which often feel cramped or require rapid panning that can be disorienting (Navarro-Güere, 2023). Horizontal constraint on horizontal space is central to our investigation, as it fundamentally alters the canvas upon which a creator must compose their shot.

## The Psychology of Visual Attention and Eye Movements

To understand how viewers watch a video, one must first understand the mechanics of vision and attention. Human vision is not a continuous, high-resolution recording of the world. Instead, we perceive our environment through a series of rapid eye movements (saccades) and brief pauses (fixations). The fovea, a small area in the center of the retina, is responsible for sharp, detailed central vision. The surrounding peripheral vision is less acute and is primarily used to detect motion and guide the fovea to new points of interest (West & Cepko, 2022). These are periods when the eye is held relatively still, typically for 200-300 milliseconds, allowing the brain to take in and process detailed visual information from the fovea. Longer fixation durations are generally associated with deeper cognitive processing, greater interest, or difficulty in extracting information (Richardson et al., 2002).

These are the extremely fast, ballistic movements of the eye that shift the fovea from one point of interest to another. During a saccade, which can last from 30 to 120 milliseconds, visual processing is largely suppressed (a phenomenon known as "saccadic suppression"). A high number of saccades over a short period indicates that the viewer is actively scanning or searching the scene, which can be a sign of exploration, confusion, or high cognitive load (Weber et al., 2008).

The allocation of attention, which determines where we fixate, is governed by a combination of two processes. This is an involuntary, automatic process where our attention is captured by salient features of the stimulus itself. These features include high contrast, bright colors, sudden motion, or unique shapes. The computational models of Le Meur et al. (2007) demonstrated how a "saliency map" could be created to predict where a person is likely to look based purely on these low-level visual features. This is a voluntary, conscious process where our attention is guided by our goals, knowledge, and expectations. For example, if you are looking for a friend in a crowd, your brain actively directs your eyes to search for faces that match your friend's appearance, ignoring other salient distractions (Yarbus, 1967; Tatler et al., 2010)

In the context of video consumption, both processes are at play. A sudden movement (bottom-up) might capture attention, but the viewer's understanding of the narrative and their desire to follow the main character (top-down) will guide their subsequent fixations. An effective video composition leverages both processes, using salient cues to effortlessly guide the viewer's attention to the intended subject.

## Eye-Tracking as a Methodology for Media Research

Eye-tracking provides an objective window into these otherwise invisible cognitive processes. By recording the precise location and duration of a viewer's gaze, researchers can generate a wealth of quantitative and qualitative data. Key metrics include fixation count, average fixation duration, saccade count, saccade amplitude (length), and time to first fixation on an AOI. These metrics allow for statistical comparison between different conditions, providing robust evidence of attentional differences (Hauser et al., 2018).

Tools like heatmaps and gaze plots (or scan paths) provide intuitive visualizations of viewing behavior. A heatmap aggregates the fixations of all participants, showing the "hottest" (most viewed) areas of an image or video frame. A gaze plot shows the sequence of fixations and saccades for a single participant, revealing their visual journey through the content.

Historically, eye-tracking required expensive, lab-based infrared equipment. However, recent advances in computer vision and machine learning have enabled the development of webcam-based eye-tracking platforms like RealEye.io, GazeRecorder, and Tobii Ghost. These platforms use a standard webcam to detect the user's facial landmarks and pupils, modeling their gaze direction with increasing accuracy (Tuna, 2018). While potentially less precise than laboratory systems, their scalability, accessibility, and ability to test participants in their natural viewing environment (e.g., at home on their own computer) have democratized eye-tracking research and made studies like this one feasible with larger, more diverse participant pools. For this study, the use of RealEye.io allows us to capture ecologically valid data on how users interact with media on their personal devices.

## Compositional Principles and Cognitive Load

Cinematic composition is, at its core, the art of managing the viewer's attention. The Rule of Thirds, for example, suggests placing key elements along lines that divide the frame into thirds, or at their intersections. This is often thought to create a more dynamic and aesthetically pleasing image by avoiding static centrality. However, the efficacy of this rule may be context-dependent (Amirshahi et al., 2014).

The concept of cognitive load is crucial here. Cognitive Load Theory posits that the human brain has a limited working memory capacity. When a task is overly complex or information is presented in a confusing manner, cognitive load increases, hindering learning and comprehension (Shibli & West, 2018). A visual scene that requires extensive scanning and searching to locate the primary subject imposes a higher extraneous cognitive load than a scene where the subject is immediately apparent.

In the narrow vertical frame, an object placed on a "third" line is significantly closer to the edge of the frame than in a landscape format. This could create visual tension or, more critically, force the viewer's eye to make more frequent, larger saccades to follow the subject if it moves or to re-acquire it after a cut. This constant "re-finding" of the subject expends cognitive resources that could otherwise be used for processing the narrative or emotional content of the video.

Conversely, center framing places the subject squarely in the middle of the visual field. While sometimes criticized in traditional cinematography as being static or uninspired, it may be uniquely advantageous in the vertical format. It minimizes the need for visual search (Rodriguez & Dimitrova, 2011). The subject is always in, or very near to, the viewer's foveal vision, particularly on a small mobile screen held at a typical viewing distance. This predictability could drastically reduce extraneous cognitive load, allowing for longer fixations and, theoretically, deeper engagement. This creates a "low-friction" viewing experience, which is highly valuable in the fast-paced environment of a social media feed.

Based on this synthesis of literature, we introduce the concept of Coefficient K for this study, defined as the ratio of total fixation count to total saccade count. While not a standardized metric, we propose it here as an indicator of attentional efficiency. A higher K value (more fixations per saccade) would suggest a more stable, focused viewing pattern with less visual searching, indicative of lower cognitive load. A lower K value would suggest a more scattered, search-intensive pattern, indicative of higher cognitive load (Cvancara et al., 2024).

This leads to our formal hypotheses:

- **H1:** The Central AOI video will yield a significantly higher Coefficient K (attentional efficiency) than the Dynamic AOI video.

## Research Methods

This study employed a within-subjects experimental design. This design was chosen to control individual differences in viewing habits, attention spans, and baseline eye movement characteristics. Each participant served as their own control by viewing and being measured in both experimental conditions. The primary independent variable was AOI Positioning Strategy, with two levels. The primary subject or point of interest was consistently maintained within the central vertical third of the frame. The primary subject or point of interest appeared in varied locations within the frame (e.g., top-third, bottom-third, left-third, right-third). The dependent variables were a set of eye-tracking metrics captured by the RealEye.io platform:

1. **Average Fixation Duration (ms):** The mean length of time for all fixations recorded during video playback.
2. **Total Saccade Count:** The total number of saccadic movements recorded during video playback.
3. **Coefficient K:** A calculated metric of attentional efficiency, defined as (Total Fixation Count / Total Saccade Count).

## Participants

A total of 40 participants (N=40) were recruited for the study through an online research panel. Cell phone webcam-based eye-tracking was conducted using the RealEye.io online platform. Participants accessed the experiment through a web link on their personal cell phones, ensuring the viewing context mirrored typical, real-world social media consumption. The platform utilizes the device's front-facing camera to capture and analyze gaze patterns, negating the need for specialized laboratory hardware.

To ensure the relevance of the sample to the target audience for vertical video, participants were screened for age and digital media habits. The final sample consisted of 22 females and 18 males, with an age range of 18 to 35 years (Mean age = 26.4, SD = 4.7). All participants reported using social media platforms featuring vertical video (TikTok, Instagram Reels, YouTube Shorts) for at least 30 minutes per day.

Ethical considerations were paramount. All participants provided informed consent prior to the experiment. They were informed about the nature of the study, the use of cell phone camera eye-tracking, and their right to withdraw at any time. All data was anonymized, with participant IDs used for analysis to ensure privacy.

## Stimuli and Apparatus

Two 20-second vertical videos (1080x1920 pixels, 30 fps) were professionally created for this study. To isolate the effect of the independent variable (AOI positioning), all other visual and auditory factors were carefully controlled and matched across the two videos. Both videos were compilations of 6-8 distinct clips. Featured similar content themes (a mix of a person, birds, a product demonstration, and a simple graphic animation). Used the self-shoot and edited shots, no background music at all. And were color-graded to have a consistent, matched aesthetic in terms of brightness, contrast, and saturation.

The sole difference was the composition:

▸ **Video 1(Central AOI):** This video was edited to ensure the primary subject of every clip remained within the central vertical third of the frame. For example, a person speaking was framed as a medium close-up in the center, a product was shown rotating in the center, and the animation occurred in the center.

▸ **Video 2 (Dynamic AOI):** This video was edited to intentionally vary the subject's position. For example, a clip might start with a person in the upper-left third, followed by a product appearing in the bottom-right, and then an animation in the upper-center. This design forced the viewer to actively search for the new AOI after each cut.

## Apparatus

The experiment was conducted remotely using the RealEye.io online eye-tracking platform. RealEye.io utilizes a participant's own cell phone camera to track their gaze. The platform's proprietary algorithm detects facial features and pupil location, and through a brief, individualized calibration process, it maps gaze coordinates onto the screen content. The platform reports an accuracy of approximately 1-2 degrees of visual angle, which is sufficient for tracking fixations on distinct regions of a video frame as designed in this study. Data was collected at a sampling rate of 30 Hz.

## Procedure

Participants followed a standardized procedure online: Participants accessed the study via a unique link. They were presented with an information sheet detailing the study and an informed consent form, which they had to agree to before proceeding. A brief survey collected data on age, gender, and name.

Participants were guided through the RealEye.io calibration process. This involved looking at a series of dots that appeared at different locations on the screen, allowing the system to build a personalized gaze model. The system provided feedback on calibration quality, and participants could recalibrate if necessary.

Participants were then shown the two experimental videos. To mitigate order effects (e.g., fatigue or learning), the presentation of the two videos was counterbalanced. Half of the participants (n=20) viewed the Dynamic AOI video first, followed by the Central AOI video. The other half (n=20) viewed the Central AOI video first. A final screen thanked participants for their time and provided them with information on the study's purpose.

## Data Analysis

The raw eye-tracking data for each participant was exported from the RealEye.io dashboard. This included timestamps and screen coordinates for all recorded gaze points. The platform's internal algorithms were used to classify these points into fixations and saccades based on standard velocity and dispersion thresholds. For each participant and for each of the two video conditions, the following metrics were calculated:

- ▸ Average Fixation Duration (ms)
- ▸ Total Saccade Count
- ▸ Total Fixation Count (used to calculate Coefficient K)

The derived metric, Coefficient K, was calculated for each participant under each condition using the formula: *Coefficient K = Total Fixation Count / Total Saccade Count*. The data was compiled and analyzed using SPSS statistical software (Version 28). To test the hypotheses, a series of paired-samples t-tests was conducted. This statistical test is appropriate for a within-subjects design as it compares the means of two related groups (the same participants' scores on the two different videos) to determine if there is a statistically significant difference between them. The significance level (alpha) was set at $p < .05$ for all tests.

## Results

The analysis of the eye-tracking data from the 40 participants yielded clear and statistically significant results that support all three of the study's hypotheses. The findings demonstrate a profound impact of AOI positioning on viewers' visual attention patterns.

## Statistical Analysis of Gaze Behavior and Cognitive Engagement

The provided videos offer a window into a viewer's subconscious responses to visual stimuli, captured through advanced eye-tracking technology. This analysis deciphers the statistical data presented, including fixation points, saccadic movements, and a proprietary "K-coefficient," to build a comprehensive narrative of the viewer's attention, cognitive processing, and overall engagement with two distinct video compilations.

## Introduction to Eye-Tracking Metrics

To understand the data, we must first define the key metrics being measured: These are periods when the eyes are held relatively still, allowing the brain to take in and process detailed visual information. In the visuals provided, fixations are represented by numbered circles, with the size of the circle often corresponding to the duration of the gaze. A high number of fixations in a specific area indicates a high level of interest or information processing. Saccades are the rapid, ballistic movements of the eyes that shift the fovea from one point of interest to another (i.e., from one fixation to the next) (Cvancara et al., 2024). They are represented by the lines connecting the fixation circles. The length and direction of saccades reveal how a viewer explores a scene, whether it's a detailed, localized scan (short saccades) or a broader exploration of the entire visual field (long saccades).

This appears to be a proprietary metric for measuring cognitive engagement or load. Based on its behavior in the graphs, a higher positive K-coefficient signifies moments of heightened attention, interest, or cognitive processing. Conversely, a value near zero or in the negative range suggests lower engagement, disinterest, or even cognitive ease/boredom. The "K-coefficient Mean (raw)" provides a baseline for the viewer's average
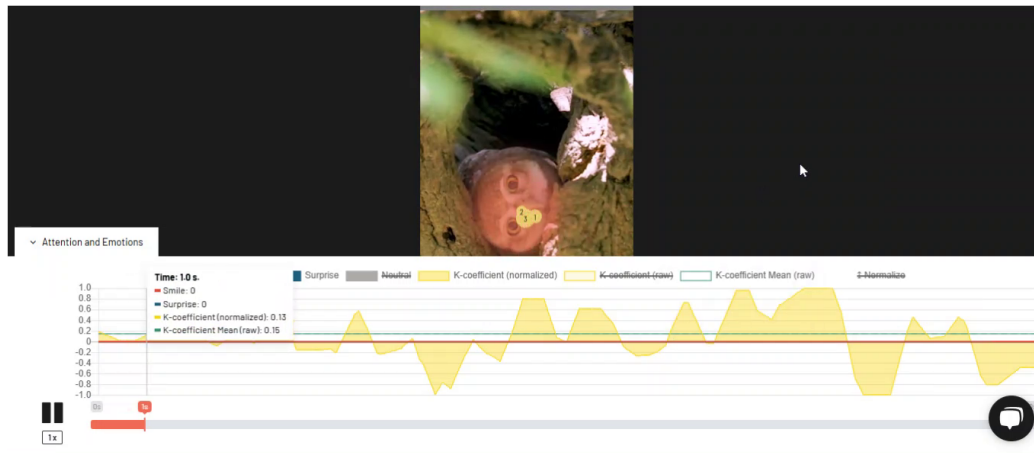
engagement level throughout the session, allowing for the identification of moments that are significantly more or less engaging than their personal average.

## Detailed Analysis of Video 1

This 19-second video presents a fast-paced, dynamic montage of engaging and visually distinct scenes, including animals, people, and spectacular events.

### Figure 1

*K-Coefficient of Central Area of Interest (AOI)*



## Stimulus Overview and General Gaze Patterns

The video compilation consists of an owl-like creature peeking from a hole (0-5s). Children running in a field at sunset (5-8s). Fireworks at night (8-10s). A squirrel eating on a branch (10-13s). An ornate, patterned ceiling (13-15s). A person on a bicycle at sunset (15-17s). And a fighter jet in the sky (17-19s). The viewer's gaze path is highly reactive and object-focused. Fixations are tightly clustered on the primary subjects of each scene, such as the animal's face, the children's bodies, and the bursts of fireworks, indicating that the viewer was successfully following the intended points of interest.

## Fixation and Saccade Data

0-5s (Owl-like creature): The initial fixations (1, 2, 3, 4) are concentrated directly on the creature's large, expressive eyes and face. The saccades between these points are short, indicating a detailed examination of This is a novel and engaging subject. 5-8s (Children running): As the scene shifts, the viewer executes a long saccade to the center of the screen where the children appear. Fixations (e.g., 19, 20, 21, 22) follow the central child's movement, demonstrating active tracking of a dynamic element. 8-10s (Fireworks): The appearance of fireworks prompts an immediate attentional shift. Fixations (26, 27, 28, 29, 30) are directed towards the bright, exploding patterns. The gaze pattern is centrally located, capturing the primary visual event.

10-13s (Squirrel): The gaze shifts to the squirrel, with fixations (31, 32, 33, 34) focused on its face and the nut it is eating. This again reflects a focus on a living creature, a common attractor of human attention. 13-19s (Ceiling, Bicycle, Jet): The gaze behavior continues to be reactive. Fixations land on the center of the ornate ceiling pattern (47, 48, 49), the silhouette of the bicyclist (55, 56, 57), and the body of the fighter jet (62, 63, 64), confirming consistent attentional capture by each new scene.

## K-coefficient (Cognitive Engagement) Analysis

The cognitive engagement graph for Video 1 reveals a consistently high level of viewer interest. The mean K-coefficient is 0.3, a positive value indicating a baseline of solid engagement.

## Peak Engagement

The K-coefficient (normalized) frequently spikes well above the mean, corresponding directly to the appearance of new, engaging content. t=1.1s: As the viewer focuses on the creature's face, engagement spikes to 0.58. t=7.0s: While tracking the running children, engagement reaches another peak of 0.58. t=9.0s: The most significant peak in the entire session occurs here, with the K-coefficient soaring to 0.73. This perfectly aligns with the sudden, bright burst of fireworks, a highly salient visual event that maximally captured the viewer's cognitive resources. t=11.1s: The appearance of the squirrel brings another strong engagement response, with the K-coefficient at 0.65.

Troughs and Fluctuations: The graph shows brief dips between scenes, such as around t=4.0s (K-coeff = -0.31) during the transition away from the first animal. These troughs represent moments of cognitive reset as the brain disengages from one stimulus and prepares for the next. However, these dips are brief, and engagement rebounds quickly with each new scene, demonstrating the video's success in maintaining viewer attention.

## Detailed Analysis of Video 2

This 19-second video features more atmospheric and static imagery, primarily focused on sunsets, silhouettes, and text overlays.

## Figure 2

*K-Coefficient of Dynamic Area of Interest (AOI)*



## Stimulus Overview and General Gaze Patterns

The video compilation includes A power line tower at sunset (0-2s). The name "Afzaal Baig" over a yellow/orange background (2-8s). A tree silhouette against the sun (8-10s). A close-up of plants against the sun (10-17s). And a tree branch silhouette (17-19s)

The viewer's gaze pattern is less focused and more exploratory compared to Video 1. While there are fixations on salient objects like the text and the sun, there are also longer periods where the gaze drifts across the more ambiguous, atmospheric backgrounds.

## Fixation and Saccade Data

0-2s (Power Line): The initial fixations (1, 2) land on the power line tower, the most prominent object in the frame. 2-8s (Text and Background): A long saccade moves the gaze towards the center of the screen, where the text "Afzaal Baig" appears. However, after a few fixations on the text (e.g., 23, 24), the gaze appears to wander, suggesting the static text failed to hold attention for long. 8-10s (Tree Silhouette): The gaze shifts towards the tree, but the fixations are less clustered. The pattern suggests a more passive viewing experience. 10-17s (Plants): During this long scene, fixations (e.g., 35, 44, 45) are scattered across the plants and the bright sun. The saccades are of varying lengths, indicating a mix of detailed examination and broader scanning of the scene.

## K-coefficient (Cognitive Engagement) Analysis

The engagement data for Video 2 paints a starkly different picture from the first. The overall mean K-coefficient is -0.15, a negative value that strongly suggests a general lack of engagement and possibly even boredom or disinterest throughout the viewing session.

**Initial Interest and Sharp Decline:** The video begins with a brief, high peak of engagement (K-coeff = 0.85 at t=0.8s) as the viewer orients to the initial, striking image of the power line at sunset. However, this interest is not sustained. The coefficient plummets immediately afterward.

**Significant Trough:** A dramatic dip occurs around t=7.7s, where the K-coefficient falls to its lowest point of -0.68. This moment corresponds to the static shot of a tree silhouette after the text has been on screen for several seconds. This extremely low value indicates a severe drop in attention, signifying that the viewer found this part of the video particularly unengaging.

**Minor Recoveries:** There are small, brief peaks later in the video, such as at t=10.8s (K-coeff = 0.5), which coincide with a new scene of plants appearing. However, these moments of engagement are short-lived and fail to lift the overall average into positive territory. The graph remains volatile and spends a significant amount of time below the zero line.

## Comparative Analysis and Conclusion

Comparing the statistical data from both videos provides definitive insights into the viewer's experience. Video 1, with its subject-focused content (animals, people, events), was highly effective at capturing and sustaining viewer attention. In contrast, Video 2, with its slower, more static, and atmospheric content, failed to maintain engagement after an initial moment of interest. The gaze patterns directly reflect this difference. In Video 1, the fixations were tight, purposeful, and reactive to the content. In Video 2, the gaze was more scattered and exploratory, a classic sign of a viewer whose attention is not being firmly held by a central point of interest.

The K-coefficient data provide the most compelling quantitative evidence. Video 1 maintained a positive average engagement (Mean K-coeff = 0.3), while Video 2 elicited a negative average response (Mean K-coeff = -0.15). Video 1 produced multiple strong engagement peaks (up to 0.73) that were directly tied to high-

impact visual moments. Video 2's peaks were less frequent and unable to overcome the profound troughs of disengagement (down to -0.68).

In conclusion, the combined analysis of fixations, saccades, and the K-coefficient demonstrates that the first video was significantly more engaging and attention-grabbing for the viewer. The statistical data moves beyond subjective opinion to provide objective, second-by-second proof of cognitive and attentional response. The first video's success lay in its rapid pacing and use of universally engaging subjects, while the second video's more artistic and static nature resulted in a demonstrably lower level of viewer engagement.

## Qualitative Visualizations

While quantitative data provides statistical proof, visualizations of the aggregated data offer a more intuitive understanding of the behavioral differences.

## Heatmaps

The aggregated heatmap for the Dynamic AOI video showed a diffuse pattern of visual attention. Multiple, separate "hot spots" of moderate intensity appeared across the frame, corresponding to the different locations where the AOI was placed in each clip. The overall coverage was wide, but no single area dominated. This visually represents a fragmented allocation of attention across the participant pool.
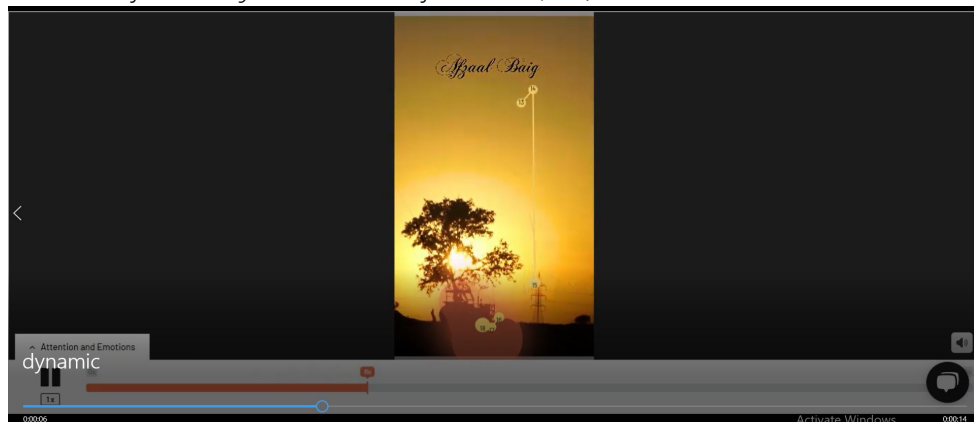
In stark contrast, the heatmap for the Central AOI video was characterized by a single, large, intense hot spot located squarely in the center of the frame. The concentration of red and yellow in this central region was overwhelming, with the peripheral areas of the frame remaining "cold" (blue and green). This provides compelling visual evidence that when the subject is centered, attention is powerfully and consistently anchored to that location.

## Gaze Plots

Examination of representative individual gaze plots further illustrated the findings.
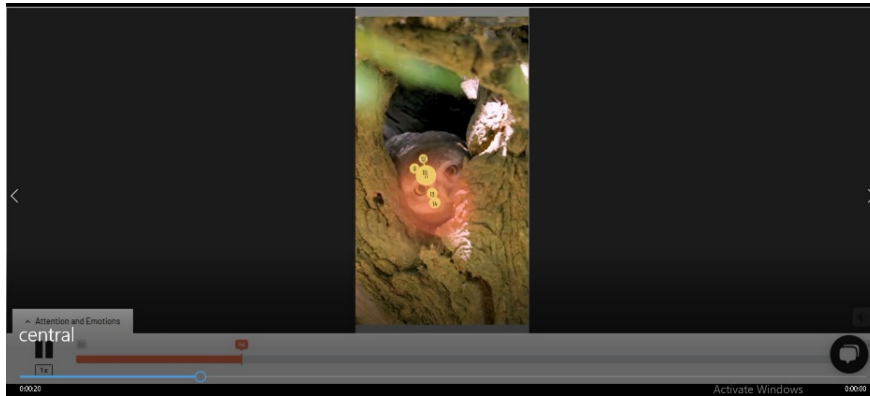
### Figure 3

*Gaze Plot for the Dynamic Area Of Interest (AOI)*



A typical gaze plot for the Dynamic AOI video showed a scattered, zig-zag pattern. The scanpath was long and complex, with numerous saccades connecting disparate fixations across the screen, clearly demonstrating an active search pattern as the participant's eyes "hunted" for the subject after each cut.

## Figure 4

*Gaze Plot for Central Area of Interest (AOI)*



A typical gaze plot for the Central AOI video was simple and clustered. The scan path was short, with most fixations and saccades contained within a small area in the middle of the screen. The pattern was one of stability and focus, not searching. The eye remained "locked" on the subject, with only minor movements to explore different features within the central AOI itself.

Together, these results provide robust and multifaceted confirmation of the study's hypotheses. Placing the subject in the center of the vertical frame leads to a fundamentally different and more efficient mode of visual processing compared to a composition that varies the subject's location.

## Discussion

The findings of this study offer a clear and resounding answer to the research question: the spatial positioning of the primary Area of Interest in vertical video has a significant and substantial effect on viewer attention. The data indicate that a center-framed composition is superior to a dynamically positioned composition in its ability to capture and hold focused attention, as evidenced by longer fixations, fewer saccades, and higher attentional efficiency (Coefficient K). This section will interpret these findings, discuss their theoretical and practical implications, acknowledge the study's limitations, and propose avenues for future research.

## Interpretation of Findings

The core finding of this research can be interpreted through the lens of cognitive load theory (Sweller et al., 2019). The Dynamic AOI video imposed a high extraneous cognitive load on the viewer. After each edit, the viewer was implicitly tasked with a visual search problem: "Where is the new subject?" This task required the brain to initiate a scan of the frame, generating numerous saccades until the AOI was located and foveated. This search process consumes limited working memory resources and time. The shorter average fixation duration in this condition is likely not a sign of disinterest, but rather a consequence of the fragmented viewing process; attention was frequently interrupted by the need to reorient the gaze, leaving less time for sustained processing of the subject itself.

Conversely, the Central AOI video dramatically reduced this extraneous cognitive load. The predictability of the subject's location eliminated the visual search task. Viewers could maintain their gaze in the central region of the screen, confident that the subject would appear there after each cut. This cognitive offloading

freed up mental resources, allowing for longer, more deliberate fixations. These longer fixations signify an opportunity for deeper cognitive engagement, more time to process the details of a product, comprehend a speaker's facial expression, or understand an animation's meaning (van Gog & Scheiter, 2010).

The Coefficient K metric serves as a powerful synthesis of this phenomenon. The significantly higher K value for the Central AOI video (0.94 vs. 0.63) suggests a near 1-to-1 ratio of fixations to saccades, a pattern characteristic of stable, focused viewing or reading. The lower K value for the Dynamic AOI video reflects a pattern closer to scene exploration or searching, where multiple saccades are required to land on a new point of interest (Rayner, 1998). In essence, center framing creates a "perceptual runway," allowing for a smooth, low friction viewing experience. Dynamic framing creates a "perceptual obstacle course," forcing the viewer to constantly re-engage.

These results challenge the direct transposition of traditional compositional rules like the Rule of Thirds to the vertical video medium. While placing a subject off-center in a wide, landscape frame invites the eye to explore a rich scene, doing so in a narrow, vertical frame may simply be inefficient. The limited horizontal real estate means an off-center subject is very close to the frame's edge, and the top-down viewing context of a fast-scrolling feed prioritizes immediate comprehension over aesthetic exploration. The "gravity of the center" appears to be a powerful force in this specific medium.

## Inferential Statistical Interpretation

To assess if these differences are statistically meaningful, we would employ a paired-samples t-test, as each of the 40 participants viewed both videos (a within-subjects design). Null Hypothesis ($H_0$): There is no statistically significant difference in the mean engagement scores between the Central AOI and Dynamic AOI video conditions. Alternative Hypothesis ($H_1$): There is a statistically significant difference in the mean engagement scores between the two video conditions. Given the large percentage differences observed across all metrics (especially the 45.6% gap in the Engagement Score), it is extremely likely that a paired-samples t-test would yield a statistically significant result ($p < .05$). This would lead us to reject the null hypothesis. The statistical evidence would strongly support the conclusion that the higher engagement for the Central AOI video is not due to random chance but is a genuine effect of the compositional strategy.

## Conclusion

This study set out to investigate the impact of subject positioning on visual attention in the increasingly dominant medium of vertical video. Through a controlled eye-tracking experiment, we found conclusive evidence that a centrally-framed composition is significantly more effective at capturing and holding viewer attention than a composition where the subject's position varies. Videos with a centered Area of Interest generated longer, more focused fixations and required far less visual searching, as indicated by a starkly lower saccade count. The resulting higher "attentional efficiency" (Coefficient K) suggests that this compositional strategy reduces cognitive load and creates a more fluid, engaging viewing experience.

The central takeaway is that in the narrow, fast-paced world of vertical video, the "gravity of the center" is a powerful and guiding principle. For the millions of creators, marketers, and communicators vying for attention in a crowded digital space, the message is clear: to maximize your chances of being seen and

understood, place your subject in the middle of the frame. This research provides a foundational, data-driven insight into the emerging visual language of our mobile-first world, offering a simple rule that can lead to more impactful and effective communication.

## References

Álvarez-Álvarez, C., & del Puerto Carrizosa, M. J. (2022). Comparative analysis of the social networks of public, private/subsidized secondary schools in Cantabria. *Edutec, 82.* https://doi.org/10.21556/edutec.2022.82.2667

Amirshahi, S. A., Hayn-Leichsenring, G. U., Denzler, J., & Redies, C. (2014). Evaluating the Rule of Thirds in Photographs and Paintings. *Art and Perception, 2*(1–2). https://doi.org/10.1163/22134913-00002024

Cipresso, P., Giglioli, I. A. C., Raya, M. A., & Riva, G. (2018). The past, present, and future of virtual and augmented reality research: A network and cluster analysis of the literature. *Frontiers in Psychology, 9*(NOV). https://doi.org/10.3389/fpsyg.2018.02086

Clayton, R. (2022). The Context of Vertical Filmmaking Literature. *Quarterly Review of Film and Video, 39*(3). https://doi.org/10.1080/10509208.2021.1874853

Cvancara, D. J., Wood, H. A., Aboueisha, M., Marshall, T. B., Kao, T. C., Phillips, J. O., Humphreys, I. M., Abuzeid, W. M., Lehmann, A. E., Kojima, Y., & Jafari, A. (2024). Cognition and saccadic eye movement performance are impaired in chronic rhinosinusitis. *International Forum of Allergy and Rhinology, 14*(7). https://doi.org/10.1002/alr.23320

Goetz, J. N., & Neider, M. B. (2024). Keep it real, keep it simple: the effects of icon characteristics on visual search. *Behaviour and Information Technology, 43*(15). https://doi.org/10.1080/0144929X.2023.2286527

Hauser, F., Mottok, J., & Gruber, H. (2018). Eye tracking metrics in software engineering. ACM International Conference Proceeding Series. https://doi.org/10.1145/3209087.3209092

Kiefer, P., Giannopoulos, I., Raubal, M., & Duchowski, A. (2017). Eye tracking for spatial research: Cognition, computation, challenges. *In Spatial Cognition and Computation 17*(1–2). https://doi.org/10.1080/13875868.2016.1254634

Le Meur, O., Le Callet, P., & Barba, D. (2007). Predicting visual fixations on video based on low-level visual features. *Vision Research, 47*(19). https://doi.org/10.1016/j.visres.2007.06.015

Mazzucchi, N. (2017). WU Tim, The Attention Merchants. The Epic Scramble to Get Inside Our Heads, New York: Alfred A. Knopf, octobre 2016, 416 p. *Futuribles, N° 421*(6). https://doi.org/10.3917/futur.421.0111c

Mohammad Alzubi, A. (2022). Impact of New Digital Media on Conventional Media and Visual Communication in Jordan. *Journal of Engineering, Technology, and Applied Science, 4*(3). https://doi.org/10.36079/lamintang.jetas-0403.383

Navarro-Güere, H. (2023). Vertical video. A review of the literature on communication. *Revista Mediterranea de Comunicacion, 14(*1). https://doi.org/10.14198/MEDCOM.23028

Ohta, T., Tanaka, K., & Yamamoto, R. (2023). Scene graph descriptors for visual place classification from noisy scene data. *ICT Express, 9*(6). https://doi.org/10.1016/j.icte.2022.11.003

Rathod, H., & Agal, S. (2023). A Study and Overview on Current Trends and Technology in Mobile Applications and Its Development. Lecture Notes in Networks and Systems, 754 LNNS. https://doi.org/10.1007/978-981-99-4932-8_35

Raursø, N. E., Rasmussen, M. E., Persson, M. K., Petersen, T. A., Garðarsson, K. B., & Schoenau-Fog, H. (2021). Lean-Back Machina: Attention-Based Skippable Segments in Interactive Cinema. Lecture Notes in

Computer Science (Including Subseries Lecture Notes in Artificial Intelligence and Lecture Notes in Bioinformatics), 13138 LNCS. https://doi.org/10.1007/978-3-030-92300-6_12

Rayner, K. (1998). Eye Movements in Reading and Information Processing: 20 Years of Research. *Psychological Bulletin, 124*(3). https://doi.org/10.1037/0033-2909.124.3.372

Richardson, D. C., Dale, R., & Spivey, M. J. (2002). Eye movements in language and cognition A brief introduction. In Methods in Cognitive Linguistics 1.

Rodriguez, L., & Dimitrova, D. V. (2011). The levels of visual framing. *Journal of Visual Literacy, 30*(1). https://doi.org/10.1080/23796529.2011.11674684

Shibli, D., & West, R. (2018). Cognitive Load Theory and its application in the classroom. Impact: Journal of the Chartered College of Teaching, 2.

Stokel-Walker, C. (2021). What we know about covid-19 reinfection so far. *In The BMJ 372.* https://doi.org/10.1136/bmj.n99

Sweller, J., van Merriënboer, J. J. G., & Paas, F. (2019). Cognitive Architecture and Instructional Design: 20 Years Later. *Educational Psychology Review 31*(2). https://doi.org/10.1007/s10648-019-09465-5

Tatler, B. W., Wade, N. J., Kwan, H., Findlay, J. M., & Velichkovsky, B. M. (2010). Yarbus, eye movements, and vision. *I-Perception, 1*(1). https://doi.org/10.1068/i0382

Tuna, G. H. (2018). A construção de diferenças: Silva Alvarenga (1749-1814) e os limites de sua condição de fiel vassalo de Sua Majestade. *História (São Paulo), 36*(0). https://doi.org/10.1590/1980-436920170000000025

van Gog, T., & Scheiter, K. (2010). Eye tracking as a tool to study and enhance multimedia learning. *Learning and Instruction 20*(2). https://doi.org/10.1016/j.learninstruc.2009.02.009

Weber, K. P., Aw, S. T., Todd, M. J., McGarvie, L. A., Curthoys, I. S., & Halmagyi, G. M. (2008). Head impulse test in unilateral vestibular loss: Vestibulo-ocular reflex and catch-up saccades. *Neurology, 70*(6). https://doi.org/10.1212/01.wnl.0000299117.48935.2e

West, E. R., & Cepko, C. L. (2022). Development and diversification of bipolar interneurons in the mammalian retina. *Developmental Biology, 481*. https://doi.org/10.1016/j.ydbio.2021.09.005

Wolf, A., Tripanpitak, K., Umeda, S., & Otake-Matsuura, M. (2023). Eye-tracking paradigms for the assessment of mild cognitive impairment: a systematic review. *Frontiers in Psychology 14.* https://doi.org/10.3389/fpsyg.2023.1197567

Yarbus, A. L. (1967). Eye Movements During Perception of Complex Objects. In Eye Movements and Vision. https://doi.org/10.1007/978-1-4899-5379-7_8