# Toward Real-World Emotion Decoding: A Transformer-Based Approach Using Movie fMRI

**Bo-Gyeom Kim [1], Patrick Styll [2], Jiook Cha [1]**

Department of Psychology, Seoul National University

Institute of Computer Engineering, TU Wien (Vienna University of Technology), Vienna, Austria

## Abstract

Real-world emotion recognition arises through continuous interactions among multiple sensory cues—dynamics often missed by standard laboratory paradigms. To investigate these dynamics, we applied a Transformer-based deep-learning model (SwiFT combined with Perceiver IO) to functional MRI data from 512 youths (ages 5–21) watching a 10-minute movie. By modeling neural signals as continuous time-series, we tracked short-term (~40s, 50 TRs) changes in seven emotions (e.g., positive, fear). Longer sequence windows and explicit hemodynamic modeling (double-gamma HRF) improved decoding accuracy, highlighting the importance of extended temporal context and precise BOLD-delay modeling. The prominent contribution of the visual cortex suggests reliance on low-level visual features within rich audiovisual stimuli. These findings demonstrate that flexible sequence-to-sequence methods effectively capture the temporal dynamics of emotion recognition under realistic conditions, deepening our understanding of real-world emotional processing.

**Keywords:** emotion perception; emotion decoding; predictive coding; constructivist theory; fMRI; deep learning; naturalistic stimuli

## Introduction

Real-world emotion recognition arises from continuous interactions among multiple sensory cues, which standard laboratory paradigms often fail to capture. Conventional neuroscience methods, typically relying on averaged or segmented neural signals, also struggle with time-series complexity. More flexible approaches are thus needed to track moment-to-moment changes under naturalistic conditions.

Several theoretical frameworks highlight the need to study emotions in continuously evolving contexts. Predictive coding proposes that the brain updates its predictions based on incoming signals, making emotional states dynamic inferences shaped by internal and external factors (Friston & Kiebel, 2009). Constructivist theory (Barrett, 2013) emphasizes personal context and past experiences, leading individuals to interpret the same stimulus differently. Together, these views suggest examining emotion perception as a process developing over time, rather than through discrete, controlled snapshots.
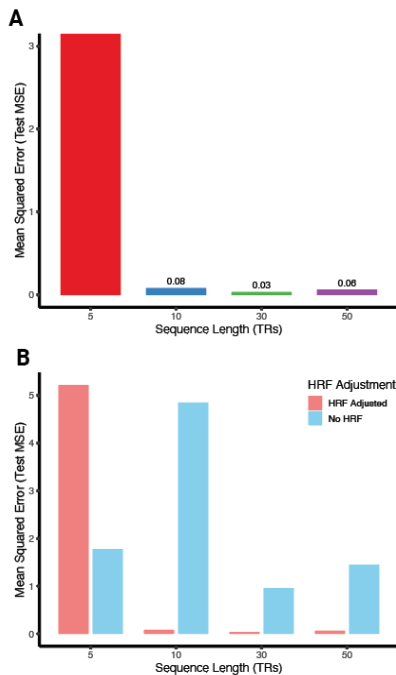
In this study, we investigate how 512 children and adolescents perceive emotion in a naturalistic setting by applying a Transformer-based model to functional MRI data collected during movie viewing. By modeling neural signals as a continuous time-series rather than averaging across discrete blocks, we aim to capture the moment-to-moment unfolding of emotion perception in an environment that more closely mirrors real-world experiences. This approach provides insight into how participants interpret evolving stimuli, bridging the gap between tightly controlled laboratory tasks and the complexities of everyday emotional encounters.

## Methods

We analyzed fMRI data from 512 youths (5–21 years, HBN dataset) including ADHD (n=301), ASD (n=63), and comorbidities (n=59). Participants watched a 10-minute *Despicable Me* (750 TRs at 0.8 s), with per-TR ratings of seven emotions obtained from human raters (Camacho et al., 2023).

To assess the feasibility of decoding emotional states, we developed a hybrid deep-learning model by combining SwiFT (Swin Transformer for fMRI; Kim et al., 2023) and Perceiver IO (Jaegle et al., 2021). SwiFT encodes spatiotemporal patches via multi-head self-attention,

84 while Perceiver IO performs sequence-to-sequence
85 regression of emotion at each TR. Multiple sequence
86 lengths (5, 10, 30, 50 TRs) and input offsets (0, 3, 5,
87 10, 20, 40 TRs) were explored to account for different
88 stimulus-response lags, and a double-gamma
89 hemodynamic response function (peak=5 s,
90 undershoot=12 s) was applied to model BOLD delays.
91 Integrated Gradients (IG) was computed for the top
92 20% of test participants (n=21) whose predictions
93 were most accurate, enabling voxel-level contribution
94 analysis for each emotion category.

95
96 Figure 1: (A) Effect of sequence length on decoding
97 performance. We tested windows of 5, 10, 30, and
98 50 TRs at a fixed learning rate, finding that longer
99 sequences consistently yielded lower MSE.
100 **(B)** Comparison of MSE between non-HRF (fixed 6 s
101 delay) and a double-gamma HRF. Modeling the full
102 hemodynamic response consistently lowered MSE,
103 highlighting the importance of accurately capturing
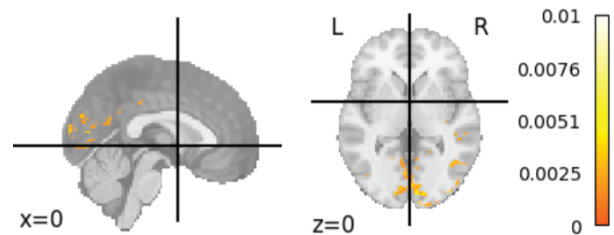104 hemodynamic delay.

## Results

106 To capture the continuous, real-world nature of
107 emotion decoding, we systematically varied
108 sequence length (5, 10, 30, 50 TRs), HRF modeling
109 (fixed 6 s delay vs. double-gamma), and time offset
110 (0–5 TRs), seeking the optimal configuration for a

111 sequence-to-sequence approach. A 50-TR window,
112 3-TR offset, and double-gamma HRF yielded the best
113 performance (MAE=0.058, MSE=0.040, r=0.894),
114 while omitting HRF correction increased error rates
115 (Figure 1).
116 Next, we assessed model performance across
117 multiple emotion dimensions. Anger ($r$ = 0748, MSE =
118 0.13), Fear ($r$ = 0.548, MSE = 0.07) showed higher
119 prediction accuracy, whereas Sad ($r$ = 0.867, MSE = 2.35),
120 Happy were more difficult to decode. These results
121 suggest that dimensions with stronger arousal cues may
122 exhibit more consistent neural signatures in this paradigm.
123 IG-based interpretation indicated that the visual
124 cortex contributed prominently to predictions across all
125 seven emotion categories, including positive and negative
126 valence. In the top-performing group of participants, this
127 region displayed consistently high voxel-wise attribution
128 scores, suggesting a broad involvement of visual
129 processing in fMRI-based emotion decoding under
130 cinematic stimulation **(Figure 2).**

131
132 Figure 2: Integrated Gradients map for Positive
133 emotion, displaying the top 5% of voxels contributing
134 to predictions. Similar patterns emerged across all
135 emotion categories, highlighting the visual cortex as
136 a key region in continuous emotion decoding.

## Discussion

138 Our findings show that longer temporal windows and
139 explicit hemodynamic modeling significantly enhance
140 continuous emotion decoding, with anger and fear
141 predicted more accurately than sad—likely reflecting the
142 stronger neural signals elicited by higher-arousal states.
143 The visual cortex consistently emerged as a key
144 contributor, aligning with prior findings on the modulation
145 of visual processing by emotion (Vuilleumier & Driver,
146 2007). Future work will compare these methods to
147 conventional approaches and investigate
148 developmental/clinical subgroups to see how different
149 populations encode emotion.

# Acknowledgements

# References

Barrett, L. F. (2013). Psychological Construction: The Darwinian Approach to the Science of Emotion. Emotion Review, 5(4), 379–389. https://doi.org/10.1177/1754073913489753

Camacho, M. C., Nielsen, A. N., Balser, D., Furtado, E., Steinberger, D. C., Fruchtman, L., Culver, J. P., Sylvester, C. M., & Barch, D. M. (2023). Large-scale encoding of emotion concepts becomes increasingly similar between individuals from childhood to adolescence. *Nature Neuroscience*, *26*(7), 1256–1266. https://doi.org/10.1038/s41593-023-01358-9

Friston, K., & Kiebel, S. (2009). Predictive coding under the free-energy principle. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *364*(1521), 1211–1221. https://doi.org/10.1098/rstb.2008.0300

Jaegle, A., Borgeaud, S., Alayrac, J.-B., Doersch, C., Ionescu, C., Ding, D., Koppula, S., Zoran, D., Brock, A., Shelhamer, E., Hénaff, O., Botvinick, M. M., Zisserman, A., Vinyals, O., & Carreira, J. (2021). Perceiver IO: A General Architecture for Structured Inputs & Outputs. *ArXiv (Cornell University)*. https://doi.org/10.48550/arxiv.2107.14795

Kim, P. Y., Kwon, J., Joo, S., Bae, S., Lee, D., Jung, Y., Yoo, S., Cha, J., & Moon, T. (2023). SwiFT: Swin 4D fMRI Transformer. *ArXiv (Cornell University)*. https://doi.org/10.48550/arxiv.2307.05916

Vuilleumier, P., & Driver, J. (2007). Modulation of visual processing by attention and emotion: windows on causal interactions between human brain regions. *Philosophical Transactions of the Royal Society B: Biological Sciences*, *362*(1481), 837–855. https://doi.org/10.1098/rstb.2007.2092