

APPENDIX PROOF FOR THEOREM 1

Author(s) Name(s)

Author Affiliation(s)

1. PROOF OF THEOREM 1

Let's consider a binary classification task as an example for simplicity. Let $(x_{A,i}, x_{V,i})$ be a point with different prediction results for audio modality A and video modality V . Assume $\exists a, b$, such that $a^T g(x_{A,i}) = -s < 0$ and $b^T h(x_{V,i}) = t > 0$ where $s, t > 0$, and the correct label is -1 . For the point-wise robustness threshold $\epsilon_{A,i}$ of this point, where an attack $\{\delta_{\epsilon_{A,i}} : \|\delta_{\epsilon_{A,i}}\|_P \leq \epsilon_{A,i}\}$ changes the prediction label. By definition, we know $a^T g(x_{A,i} + \delta_{\epsilon_{A,i}}) \geq 0$ and $a^T g(x_{A,i} + \delta) \leq 0$ for all $0 \leq \delta \leq \epsilon_{A,i}$. If $s < t$, then the fused network predicted the wrong label even without any noise.

$$f(x_{A,i}, x_{V,i}) = (a, b)^T (g(x_{A,i}) \oplus h(x_{V,i})) \quad (1)$$

$$= a^T g(x_{A,i}) + b^T h(x_{V,i}) \quad (2)$$

$$= -s + t > 0 \quad (3)$$

Otherwise, in the case of $s \geq t$, by applying Intermediate Value Theorem to $g(x)$, there exists a point $0 \leq \delta \leq \epsilon_A$ such that $a^T g(x_{A,i} + \delta) = -t/2$:

$$\begin{aligned} f(x_{A,i} + \delta, x_{V,i}) &= (a, b)^T (g(x_{A,i} + \delta) \oplus h(x_{V,i})) \\ &= a^T g(x_{A,i} + \delta) + b^T h(x_{V,i}) \\ &= -\frac{t}{2} + t > 0 \end{aligned} \quad (4)$$

In both cases, we could find a noise $0 \leq \delta < \epsilon_{A,i}$ within the original unimodal robustness threshold to attack the multimodal network successfully. Vice versa for video. Thus, a unimodal attack can break a multimodal model, which we also empirically verified the existence of such cases in our experiments. Section ??.

We postulate that such phenomenon is like the McGurk Effect, where multimodal fusion would further distort the already non-convex decision boundary (Figure ??), making the fused decision boundary very different than the original ones and unpredictable.

2. REFERENCES

Waveform complexity: A new metric for EEG analysis