

What Really Carries Emotion in TTS?

Zonos EDA, Kipling matching, and artifact analysis

Kacper Wikiel

IYPT-style acoustic research

May 2026

Research Question

Question

Does emotional TTS encode emotion mainly through pitch, tempo, energy, pauses, or timbre?

Original claim tested here

- Emotional control may not be a full actor-like performance.
- A model may simulate emotion by moving a few acoustic levers.
- We need measured audio features, not definitions or vendor claims.

Hypotheses

- ① **Simple-channel hypothesis:** intended emotion is decodable from a small subset of acoustic channels.
- ② **Nonlinearity hypothesis:** increasing an emotion vector is not a linear amplification of one feature.
- ③ **Channel mismatch hypothesis:** the channel that moves most physically is not necessarily the channel that best classifies emotion.

What would make this interesting?

If timbre changes most, but pitch alone classifies emotion best, the model is not just doing one obvious prosody trick.

Actual Experiment Run

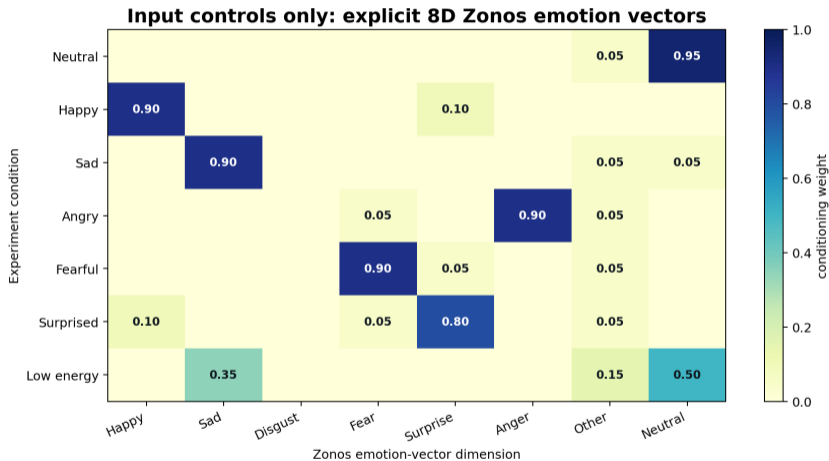
Generated dataset

- Model: Zonos v0.1 transformer
- 5 fixed English sentences
- 7 conditions: neutral, happy, sad, angry, fearful, surprised, low energy
- 35 base WAV files
- 12 additional sweep WAV files

Controls

- text held fixed
- Zonos emotion vector changed
- pitch_std fixed at 35
- speaking_rate fixed at 15
- CFG fixed at 2.0
- short-sentence EDA used Zonos' default speaker path
- later poem matching fixes speaker identity explicitly

Input Control, Not a Correlation Matrix



Rows are experiment conditions. Columns are Zonos emotion-vector dimensions. The result matrix starts later.

Feature Extraction

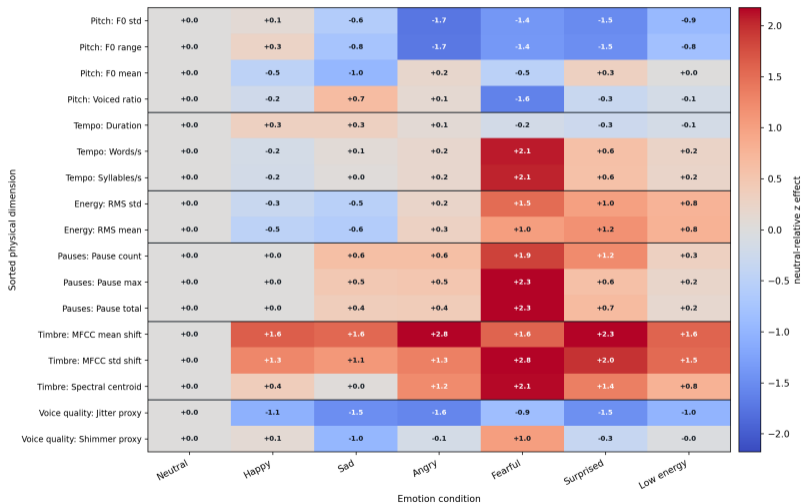
Feature map: each metric is a physical hypothesis

The analysis asks which acoustic channel moves when the model is told to synthesize an emotion.



EDA: Emotions on X, Physical Dimensions on Y

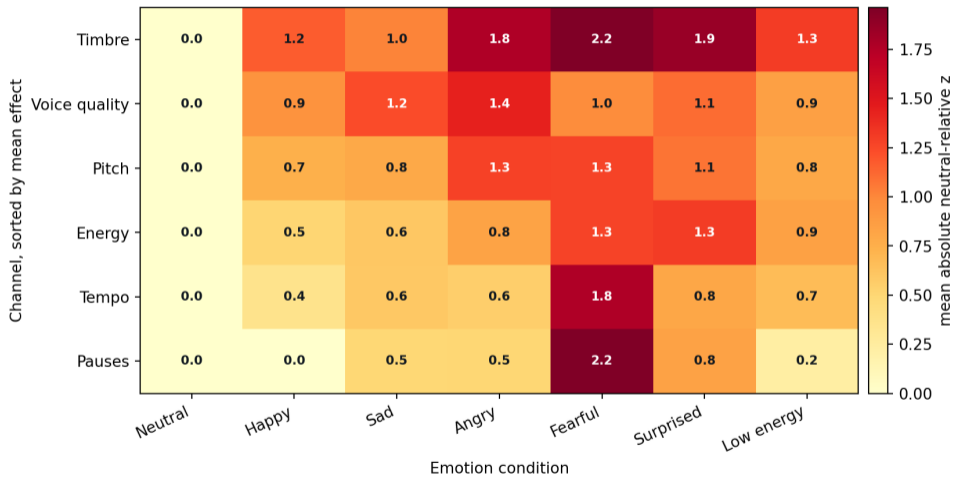
Cell value = mean normalized change relative to the neutral output of the same sentence.



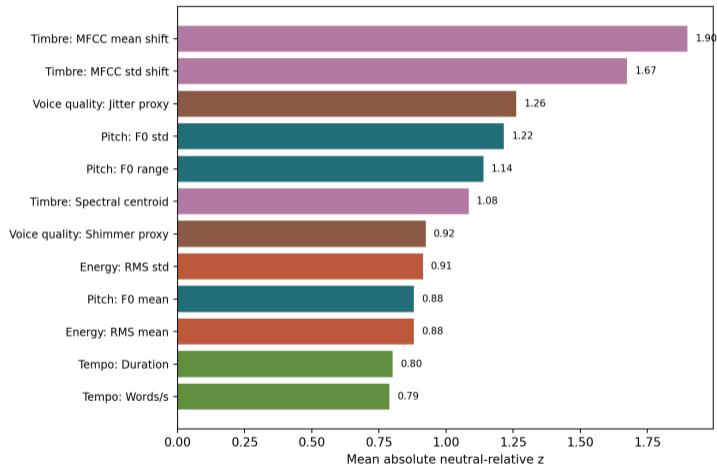
First Read of the EDA

- Fearful is the strongest global perturbation: pauses, tempo, energy, timbre, and centroid all move.
- Happy is not simply “higher pitch”: its F0 mean goes down on average, while MFCC shift is large.
- Angry, fearful, and surprised reduce F0 std/range relative to neutral in this run.
- Timbre moves for nearly every non-neutral condition, so emotion is not only speed or pitch.
- Neutral is zero by construction because every sample is compared against its same-sentence neutral baseline.

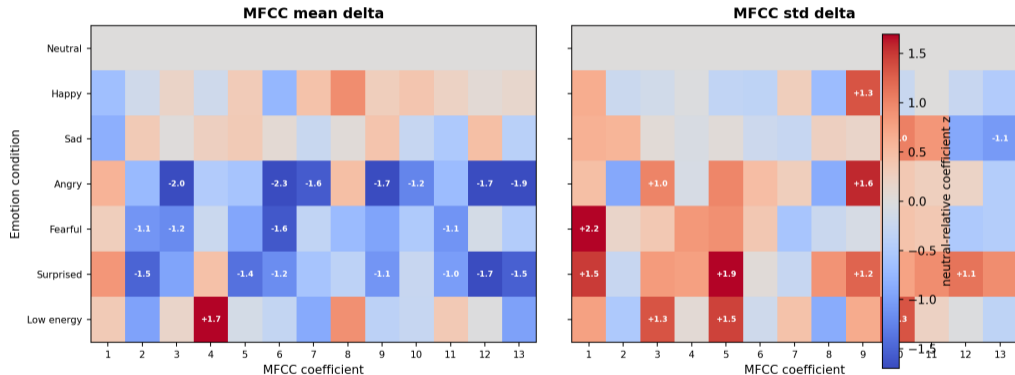
Which Channel Moves Most?



Top Moving Physical Dimensions

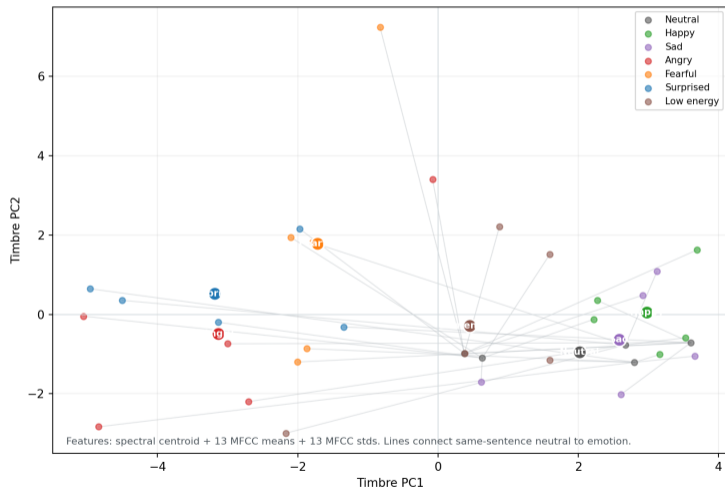


Timbre Fingerprint: MFCC Changes



MFCC coefficients are a compact spectral-envelope vector: a better timbre object than raw spectrogram pixels.

Timbre Space



How to Read Timbre Visuals

- MFCC coefficients summarize the spectral envelope, so they are closer to barwa than raw waveform energy.
- The MFCC heatmap is a fixed-size vector view: it can be averaged, tested with ANOVA, and fed to a classifier.
- The PCA plot uses only timbre features, so any separation there is not caused directly by pitch or speech rate.
- This is stronger evidence than spectrogram subtraction because the axes are stable physical descriptors.

Statistical Results

Test	Strongest features	Score	Interpretation
ANOVA	F0 std	F=9.01, p=1.62e-05	Pitch variability differs strongly between emotion conditions.
ANOVA	MFCC std shift	F=8.39, p=2.97e-05	Timbre distribution changes are significant.
ANOVA	MFCC mean shift	F=7.56, p=6.94e-05	Voice color shifts across conditions.
ANOVA	Pause total	F=5.70, p=5.70e-04	Some emotions insert timing structure.
Mutual information	F0 std, F0 range	MI=0.635, 0.582	Pitch variation carries the most label information.
Mutual information	MFCC std, MFCC mean	MI=0.521, 0.467	Timbre is the second strongest signal.

Channel Ablation Classifier

Leave-one-sentence-out 7-class classification

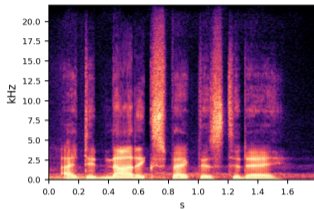
Random baseline is about 0.14. A channel is useful if it predicts emotion on held-out text.

Feature channel	Accuracy
Pitch only	0.60
Timbre only	0.37
Voice quality only	0.26
Pauses only	0.23
Energy only	0.20
Tempo only	0.17

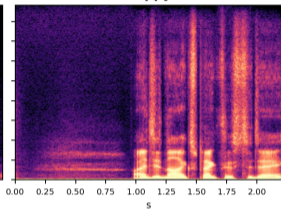
Pitch is the best emotion label decoder, even though timbre has the largest physical movement.

Spectrograms: One Fixed Sentence

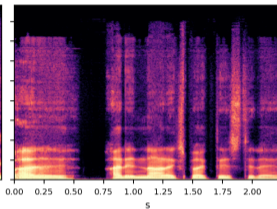
Neutral



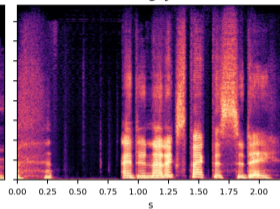
Happy



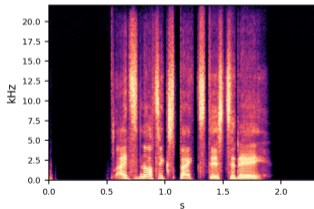
Sad



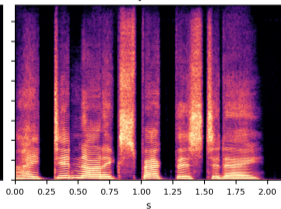
Angry



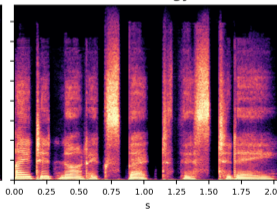
Fearful



Surprised



Low energy



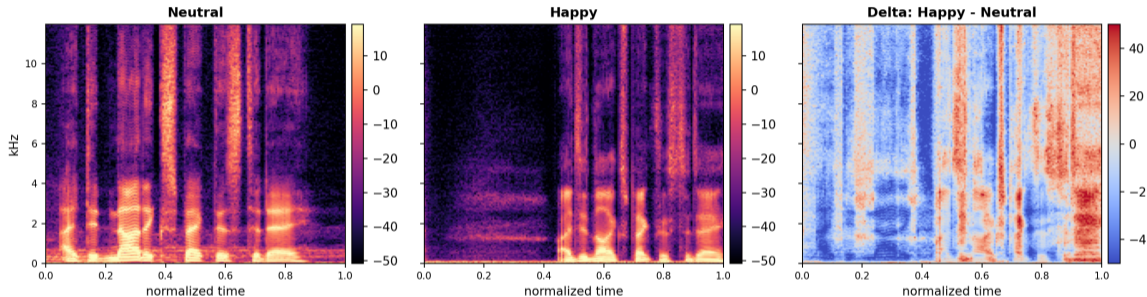
Spectrogram Observations

- Fearful contains a long silence block before speech, matching the pause spike in the heatmap.
- Angry and surprised show stronger high-frequency structure than neutral.
- Happy changes the time envelope and high-frequency texture, not just pitch.
- Low energy is not simply quieter; it still shows strong vertical energy bursts.

Interpretation

The emotion vector changes the temporal and spectral shape of the utterance. The result is richer than a pure speed knob, but still measurable through simple acoustic channels.

Delta Spectrogram: Visual Sanity Check



Naive visual subtraction after time-normalization. Useful as a sanity check, not as the main metric.

This is useful to see that something changed, but it is not the main evidence.

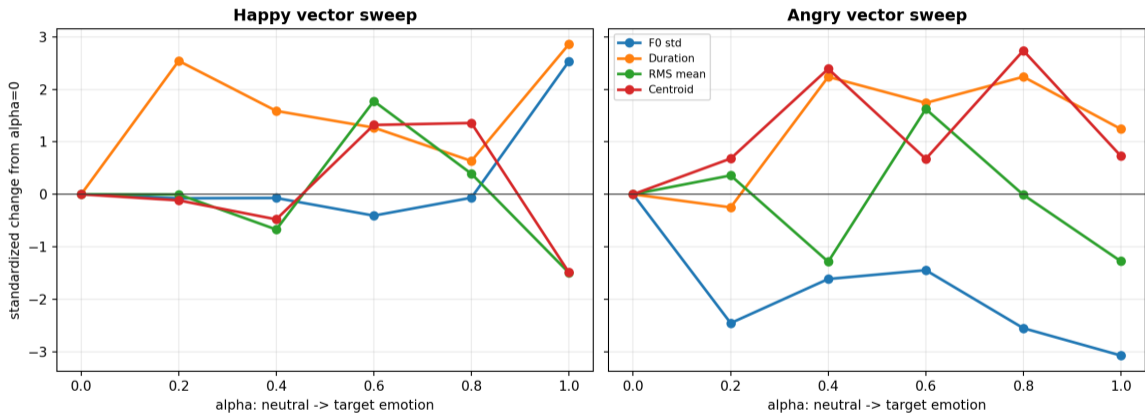
Why Spectrogram Delta Is Not the Metric

Skeptical point

Subtracting spectrogram images is visually intuitive, but it does not scale cleanly into a physical conclusion.

- Emotional TTS changes duration, pauses, and phoneme timing, so frames are not naturally aligned.
- A large delta may mean different timing, not different timbre.
- The image is high-dimensional and hard to reduce to a stable vector across many samples.
- Therefore the deck uses MFCC, spectral centroid, F0, RMS, pause, and rate features as the quantitative evidence.

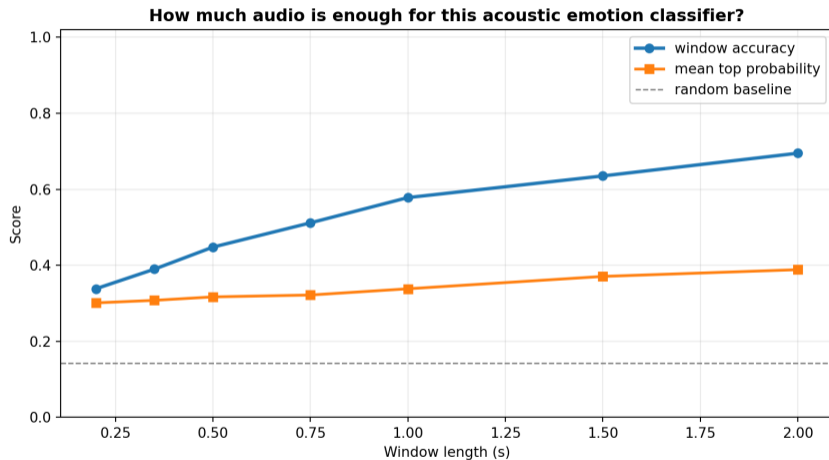
Nonlinearity Probe



Nonlinearity Result

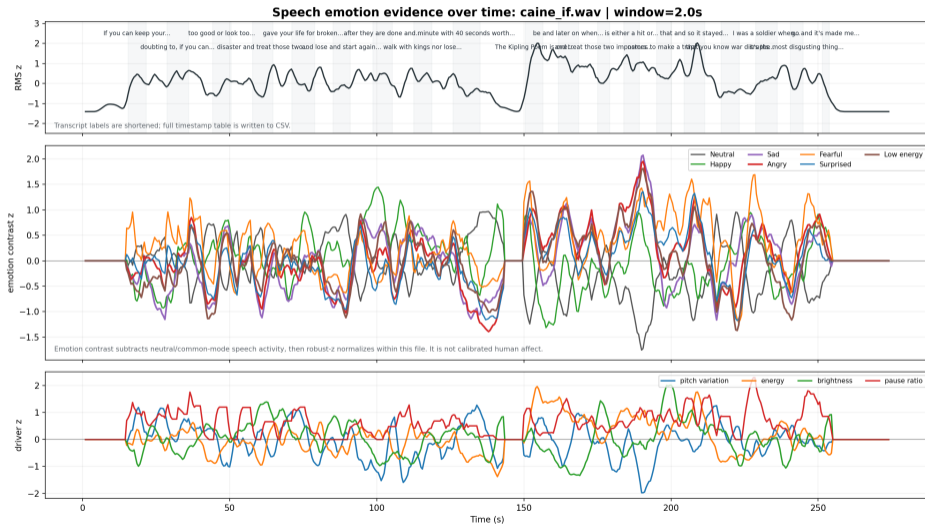
- The happy sweep is not monotonic: duration jumps early, F0 std changes mostly near the endpoint.
- The angry sweep is also not a straight line: centroid and duration peak at intermediate settings.
- Holding the random seed fixed across alpha values reduces the sampling confound.
- This supports the nonlinearity hypothesis: the emotion vector changes model state, not one scalar acoustic knob.

How Long Is Enough?

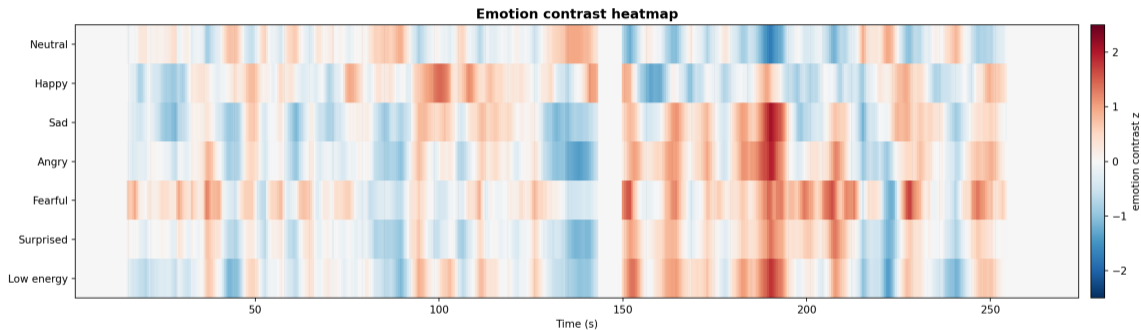


In this local acoustic classifier, 1.0 s is usable but weak; 1.5–2.0 s is more defensible for a timeline.

Real Speech Check: Emotion Over Time



Real Speech Check: Contrast Heatmap



Contrast subtracts neutral/common speech activity, so the plot shows what increases locally, not just loudness.

From EDA to Matching

New task

Use the Michael Caine recitation of Kipling's *If* as a target acoustic-emotion trajectory, then test whether controllable TTS can reproduce that trajectory.

- Target is not a single emotion label for the whole poem.
- Each canonical poem line gets a vector of evidence: emotion contrasts plus acoustic drivers.
- Dominant emotion is only a compressed summary used for quick comparison.
- The real research question becomes: can TTS follow a time-varying acoustic performance?

Caine Target: Canonical Poem Lines

Kipling 'If' recitation: canonical poem lines colored by strongest rising acoustic emotion

Stanzas 1-2

81	015.4-018.1s	Fearful +0.66	If you can keep your head when all about you
82	018.1-021.6s	Neutral +0.18	Are losing theirs and blaming it on you,
83	021.6-026.2s	Fearful +0.29	If you can trust yourself when all men doubt you,
84	026.2-029.3s	Fearful +0.45	But make allowance for their doubting too;
85	029.8-033.5s	Fearful +0.49	If you can wait and not be tired by waiting,
86	033.5-037.2s	Fearful +0.58	Or being lied about, don't deal in lies,
87	037.2-041.1s	Fearful +0.73	Or being hated, don't give way to hating,
88	041.1-045.8s	Neutral +0.18	And yet don't look too good, nor talk too wise:
89	047.2-050.4s	Happy +0.68	If you can dream—and not make dreams your master;
90	050.4-054.1s	Neutral +0.17	If you can think—and not make thoughts your aim;
91	054.1-058.1s	Fearful +0.51	If you can meet with Triumph and Disaster
92	058.1-062.1s	Neutral +0.12	And treat those two impostors just the same;
93	062.1-065.4s	Fearful +0.60	If you can bear to hear the truth you've spoken
94	065.4-069.0s	Neutral +0.05	Twisted by knaves to make a trap for fools,
95	069.0-072.6s	Neutral +0.23	Or watch the things you gave your life to, broken,
96	072.6-076.8s	Happy +0.54	And stoop and build 'em up with worn-out tools:

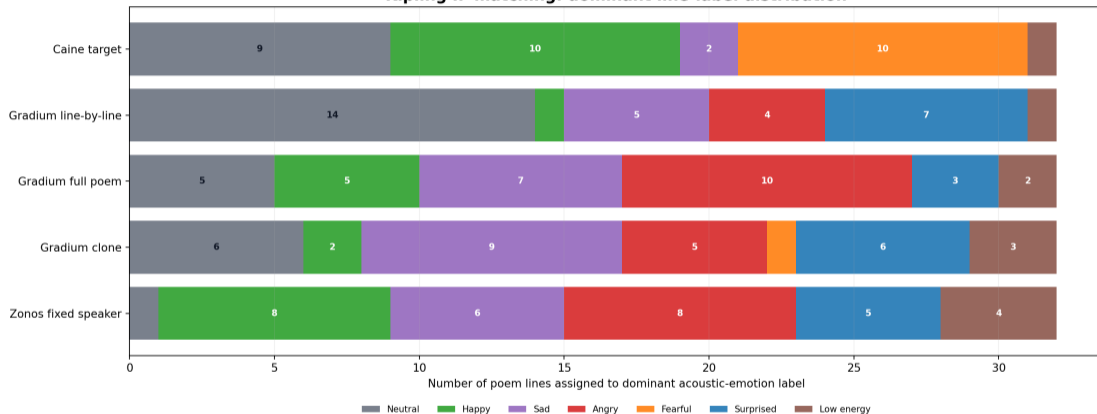
Stanzas 3-4

17	076.8-080.8s	Happy +0.30	If you can make one heap of all your winnings
18	080.8-084.1s	Neutral +0.11	And risk it on one turn of pitch-and-toss,
19	085.1-087.2s	Happy +0.26	And lose, and start again at your beginnings
20	087.6-091.0s	Neutral -0.04	And never breathe a word about your loss;
21	092.0-095.1s	Low energy +0.49	If you can force your heart and nerve and sinew
22	095.5-099.5s	Happy +1.09	To serve your turn long after they are gone,
23	099.7-103.5s	Happy +1.17	And so hold on when there is nothing in you
24	103.5-107.6s	Happy +0.57	Except the Will which says to them: "Hold on!"
25	107.6-111.9s	Happy +0.76	If you can talk with crowds and keep your virtue,
26	111.9-114.7s	Sad +0.60	Or walk with Kings—nor lose the common touch,
27	114.7-119.2s	Sad +0.49	If neither foes nor loving friends can hurt you,
28	119.2-123.0s	Happy +0.29	If all men count with you, but none too much;
29	123.0-126.2s	Fearful +0.42	If you can fill the unforgiving minute
30	126.2-130.7s	Fearful +0.43	With sixty seconds' worth of distance run,
31	130.7-135.8s	Neutral -0.09	'Tis yours is the Earth and everything that's in it,
32	135.8-142.3s	Happy +0.34	And—which is more—you'll be a Man, my son!

Color means strongest positive non-neutral emotion contrast per line; gray means no clear rise. Scores are relative acoustic evidence, not human-rated emotion.

Caine Target Distribution

Kipling IF matching: dominant line-label distribution



Caine target is mostly fearful, happy, and neutral. Most TTS attempts drift into other acoustic regions.

Gradium: Useful Negative Control

What Gradium controls

- voice choice and custom voice cloning
- speed via padding controls
- sampling temperature
- CFG / voice similarity
- text and pause tags

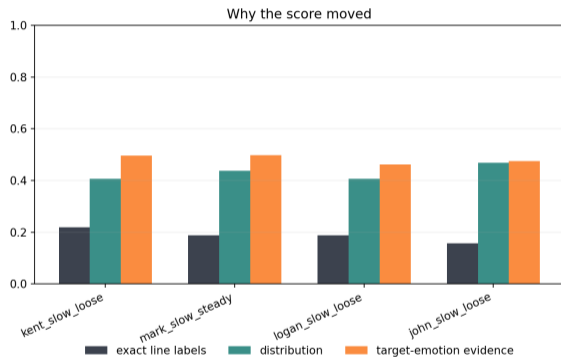
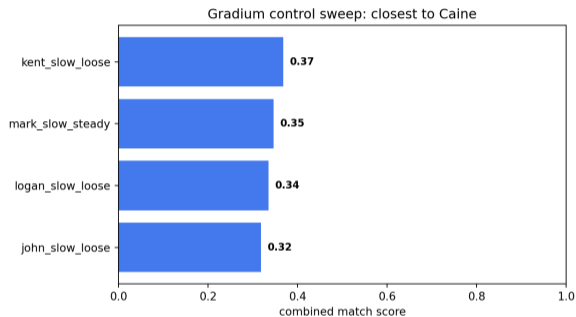
What it did not give us

- no direct emotion-vector trajectory
- cloning changed speaker/timbre, not performance arc
- full-poem generation beat line-by-line prosody
- but still did not converge to Caine's line map

Result

Gradium is a good voice/timbre baseline, but not the right primary model for causal emotion-control research.

Gradium Sweep Results



Duration was removed from the score; the comparison is emotion-pattern only.

Zonos Matching Design

Why Zonos

- explicit 8D emotion vector
- direct fear, happiness, sadness, anger, surprise, neutral axes
- pitch variation and speaking rate are separate controls

Poem generation

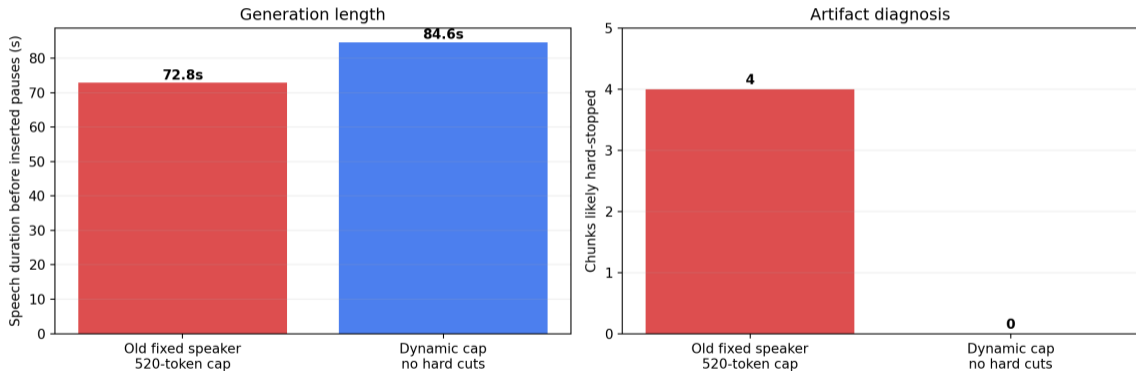
- target = Caine line labels
- consecutive same-label lines merged into chunks
- fixed Kent speaker embedding
- dynamic token cap to avoid clipped chunks

Important correction

Without a fixed speaker embedding, Zonos sounded like multiple people. Emotion matching is meaningless unless identity is controlled.

Artifact Diagnosis

Why the first Zonos poem sounded cut: `max_new_tokens` was too low



The first fixed-speaker Zonos poem had hard stops because `max_new_tokens=520` capped long chunks at about 6 seconds.

Zonos Fixed Speaker: Dominant Line View

Zonos Caine Target Runs: canonical poem lines colored by strongest rising acoustic emotion

Stanzas 1-2

81	000.0-001.8s	Low energy +0.77	If you can keep your head when all about you
82	001.8-004.6s	Sad +0.68	Are losing theirs and blaming it on you,
83	005.2-007.1s	Angry +1.13	If you can trust yourself when all men doubt you,
84	007.6-009.1s	Angry +1.18	But make allowance for their doubting too;
85	009.7-011.7s	Angry +1.44	If you can wait and not be tired by waiting,
86	011.7-015.2s	Angry +1.42	Or being lied about, don't deal in lies,
87	015.2-017.4s	Angry +1.23	Or being hated, don't give way to hating,
88	017.4-020.4s	Angry +0.43	And yet don't look too good, nor talk too wise:
89	021.3-024.1s	Happy +0.63	If you can dream—and not make dreams your master;
90	024.9-027.3s	Low energy +0.25	If you can think—and not make thoughts your aim;
91	027.7-029.6s	Surprised +0.54	If you can meet with Triumph and Disaster
92	029.6-032.4s	Surprised +0.52	And treat those two impostors just the same;
93	032.9-035.1s	Surprised +0.55	If you can bear to hear the truth you've spoken
94	035.8-037.5s	Angry +0.52	Twisted by knaves to make a trap for fools,
95	037.5-041.4s	Angry +0.47	Or watch the things you gave your life to, broken,
96	041.4-044.2s	Neutral +0.17	And stoop and build 'em up with worn-out tools:

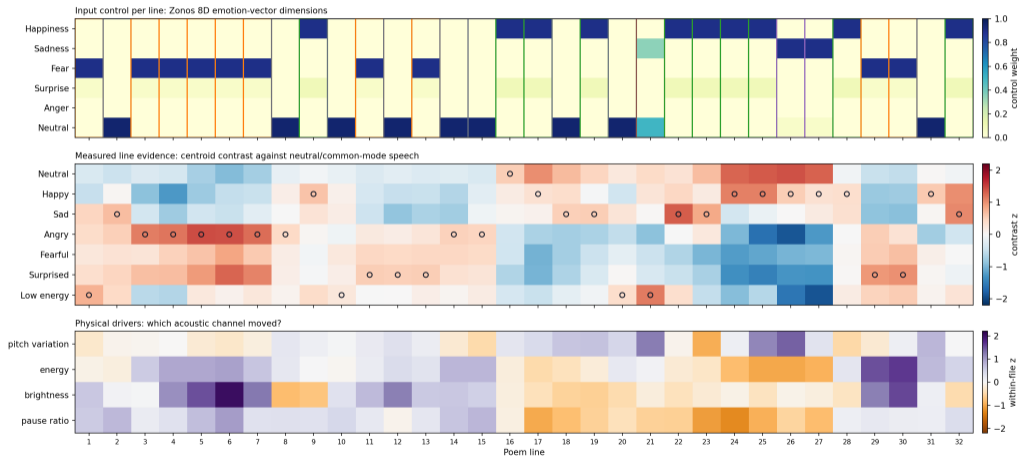
Stanzas 3-4

17	044.9-046.5s	Happy +0.34	If you can make one heap of all your winnings
18	046.7-049.3s	Sad +0.52	And risk it on one turn of pitch-and-toss,
19	050.2-051.2s	Sad +0.53	And lose, and start again at your beginnings
20	052.0-054.6s	Low energy +0.38	And never breathe a word about your loss;
21	054.6-058.3s	Low energy +1.14	If you can force your heart and nerve and sinew
22	058.8-061.1s	Sad +1.34	To serve your turn long after they are gone,
23	061.1-063.1s	Sad +0.83	And so hold on when there is nothing in you
24	063.1-066.0s	Happy +1.11	Except the Will which says to them: "Hold on!"
25	066.0-068.2s	Happy +1.11	If you can talk with crowds and keep your virtue,
26	068.2-072.9s	Happy +0.53	Or walk with Kings—nor lose the common touch,
27	073.6-076.1s	Happy +0.33	If neither foes nor loving friends can hurt you,
28	076.9-079.2s	Happy +0.37	If all men count with you, but none too much;
29	079.5-081.0s	Surprised +0.96	If you can fill the unforgiving minute
30	081.0-083.3s	Surprised +1.01	With sixty seconds' worth of distance run,
31	084.1-086.3s	Happy +0.57	'Tis yours is the Earth and everything that's in it,
32	086.3-089.0s	Sad +1.05	And—which is more—you'll be a Man, my son!

Color means strongest positive non-neutral emotion contrast per line; gray means no clear rise. Scores are relative acoustic evidence, not human-rated emotion.

Zonos Fixed Speaker: Multidimensional View

Zonos Kipling match: input vector, measured emotion evidence, and acoustic drivers



Rows 1-6 are the intended Zonos controls. Rows below are measured from audio, so disagreement reveals entanglement/artifacts rather than a label bug.

Why Dominant Emotion Is Not Enough

- A line can have rising fearful, angry, and low-energy evidence at the same time.
- Dominant label discards the second and third strongest axes.
- Physical drivers can disagree with the emotional label: pitch variation, energy, brightness, and pause ratio move separately.
- The multidimensional plot is therefore the real evidence view; the colored transcript is only the readable summary.

Research implication

TTS emotion is not a one-hot state. It is a vector-valued acoustic transformation, and artifacts can masquerade as emotion if not audited.

Matching Results So Far

Run	Duration	Exact	Rank	Dist.	Main observation
Gradium line-by-line	140.8 s	0.219	0.480	0.406	voice/settings sweep; no fear channel reproduced
Gradium full poem	134.4 s	–	–	–	better continuity, but drifts to angry/sad/happy
Gradium clone	138.7 s	0.125	0.337	0.375	cloning changes timbre, not Caine-like performance
Zonos fixed speaker, clipped	77.5 s	0.219	0.449	0.563	invalid: several chunks hard-stopped near 6 s
Zonos fixed speaker, dynamic cap	89.3 s	0.188	0.428	0.375	valid audio; still not converging to Caine

Exact = same dominant line label. Rank = target label ranked high even if not dominant. Dist. = label-distribution similarity.

Interpretation of the Matching Failure

- Explicit emotion vectors do move Zonos acoustics, but not like a human actor following the poem.
- Stronger generation hygiene made the audio more valid, but the match score dropped.
- That is not a failure of the analysis: it means previous artifacts were inflating the apparent match.
- Caine's performance contains continuity, rhetoric, breath, emphasis, and line-to-line intent, not just per-line emotion labels.

Updated thesis

Current TTS controls can steer acoustic channels, but matching a real actor requires controlling a trajectory of multiple physical dimensions, not only selecting emotion labels.

Answer to the Research Question

Measured answer

In the short-sentence Zonos EDA, emotion is carried most visibly by timbre, but most decodably by pitch variation. In the Kipling matching task, this is insufficient: real performance is a multidimensional trajectory.

- Largest physical movement: **timbre** (mean absolute neutral-relative $z = 1.55$).
- Best single-channel classifier: **pitch** (leave-one-sentence-out accuracy = 0.60).
- Strongest individual statistical feature: **F0 std** (ANOVA $p = 1.62e-05$, MI = 0.635).
- Matching result: direct emotion vectors alone do not reproduce Caine's poem trajectory.

What Is Still Missing

- Repeat with multiple speaker embeddings to separate emotion from speaker-style drift.
- Add ASR WER to test whether stronger emotion damages intelligibility.
- Add speaker embedding similarity to test whether emotional control changes identity.
- Add parameter sweeps over Zonos emotion intensity, pitch_std, speaking_rate, and CFG.
- Compare against IndexTTS-2, Chatterbox, OmniVoice, FastPitch, and FastSpeech 2.
- Add human ratings: perceived emotion, naturalness, and acting quality.

CodeSOTA Improvement

Gap

A TTS leaderboard should not only rank naturalness. It should expose controllability: emotion, pitch, tempo, energy, pauses, timbre, reference audio, and tags.

- Add structured model fields for each control surface.
- Add a small reproducible benchmark: same sentences, same controls, same feature extractor.
- Report both “what moves” and “what classifies” because those are different claims.

Takeaway

Defensible thesis

Emotional TTS in this experiment is not a full hidden theatrical performance. It is a structured acoustic transformation: strong timbre movement plus pitch features that make the intended label easy to decode.

The research value is the mismatch: the biggest physical effect is timbre, while the strongest classifier channel is pitch.

Artifacts and Sources

- Generated data: `runs/zonos_emotion_research/`
- Kipling matching: `runs/zonos_kipling_match/`
- Gradium controls/cloning: `runs/gradium_caine_match/`
- Experiment script: `scripts/zonos_emotion_research.py`
- Matching script: `scripts/zonos_kipling_match.py`
- Slides: `slides/emotion_tts_iyp.pdf`
- Zonos: <https://github.com/Zyphra/Zonos>
- Zonos conditioning: https://github.com/Zyphra/Zonos/blob/main/CONDITIONING_README.md
- OmniVoice: <https://github.com/k2-fsa/OmniVoice>
- CodeSOTA TTS: <https://www.codesota.com/text-to-speech>
- Real speech sample: <https://www.youtube.com/watch?v=EEFMVif12UY>