

Recomposer: Event-Roll-guided Generative Audio Editing

Dan Ellis *Google DeepMind*

Joint work with **Edu Fonseca, Ron Weiss, Kevin Wilson, Pascal Getreuer, Scott Wisdom, Hakan Erdogan, John Hershey, Aren Jansen, Channing Moore, Manoj Plakal**

Mila Conversational AI Reading Group, Jan 29th 2026

<https://arxiv.org/abs/2509.05256>



The Team



Dan Ellis



Eduardo Fonseca



Ron Weiss



Kevin Wilson



Pascal Getreuer



Scott Wisdom



Hakan Erdogan



John Hershey



Aren Jansen



Channing Moore



Manoj Plakal

Context: Machine Learning for Audio

Domain

- Speech
- Music
- **Environmental sounds**

Applications

- Classification / Understanding
- Generation
- **Editing / Modification**

No LLM!

Scenario: Sound Recomposition

Select clip

0022bc1c4126c0a0_30000@8008

00080086c5823584_130000@1994	!Brief tone; Speech; Music; Explosion; <u>Sneeze</u>
0022bc1c4126c0a0_30000@8008	!Generic impact sounds; Speech; Mechanisms; Motorcycle; <u>Sneeze</u> ; Traffic noise, roadway noise
00b7026f61c9986f_70000@6786	Music; Singing; <u>Sneeze</u>
020092042b0790c9_60000@6659	Stream, river; Music; <u>Sneeze</u> ; Sound effect; Liquid; Bell
020d6108c3f55879_0@4275	!Brief tone; Sound effect; Telephone; <u>Sneeze</u> ; Noise
02dae3eb9a91ac1a_200000@5556	<u>Sneeze</u> ; Bell; Train
03105a88f52883d2_380000@1078	!Generic impact sounds; Speech; Breathing; Laughter; Mechanisms; <u>Sneeze</u> ; Train
0345d7995b6cd72d_30000@5693	Music; Singing; <u>Sneeze</u>
034720c57c16de65_40000@2938	Speech; Breathing; <u>Sneeze</u> ; Mechanisms
034865f16dd11797_0@8331	Engine; Speech; !Human group actions; Steam; <u>Sneeze</u> ; Train
036dd4d2829e6c1f_30000@6551	Singing; Whistling; Music; <u>Sneeze</u> ; Hands
03e320f85c129225_30000@4292	Music; Stream, river; <u>Sneeze</u>
042f0f2cebe8f2b9_20000@4615	Speech; Music; <u>Sneeze</u>
050f3a218d0cf3ad_30000@6417	Speech; <u>Sneeze</u> ; Wind; Boat, Water vehicle
051bb169863609c6_30000@5044	Speech; Singing; Whistling; Music; <u>Sneeze</u>
0599626613d898e4_90000@5317	Speech; Laughter; Dog; Typing; <u>Sneeze</u>
0599626613d898e4_90000@6064	Speech; Laughter; Dog; Explosion; <u>Sneeze</u>
065b1b7b28a65031_30000@4657	Singing; Music; <u>Sneeze</u> ; !Human group actions; Hands

0022bc1c4126c0a0_30000@8008

0:00 / 0:10

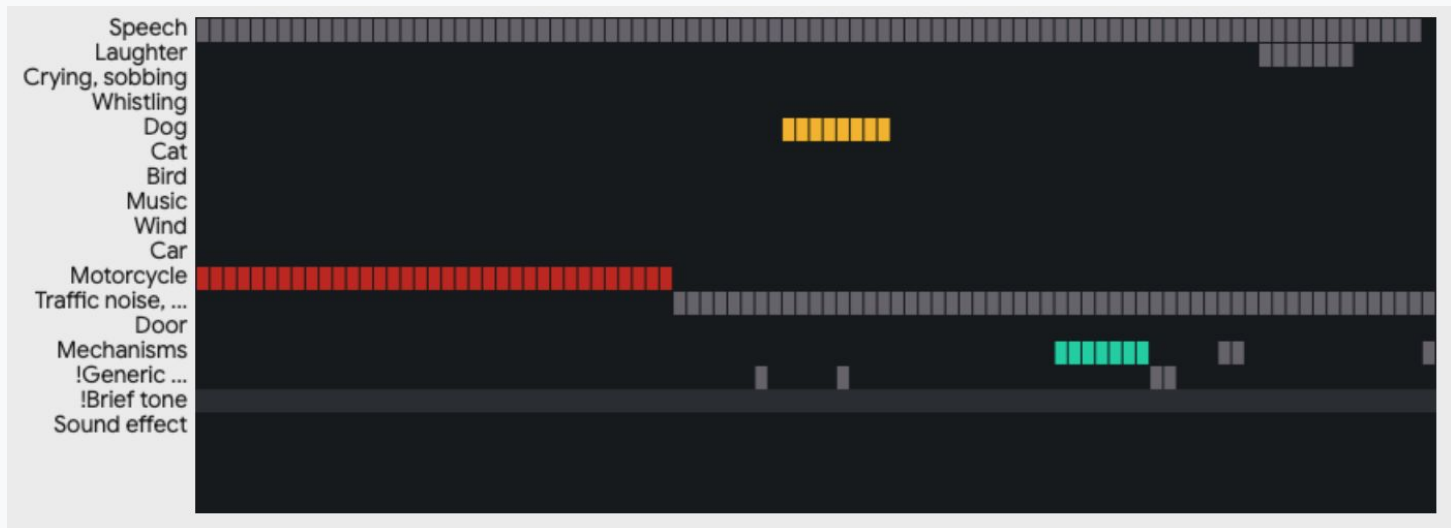
> 3.2 Run ReComPosiTron

(4) Decode parameters

This is a screen recording of the refreshed go/recompositron-demo

100% completed at 1.30 PM

Contribution: Event Roll

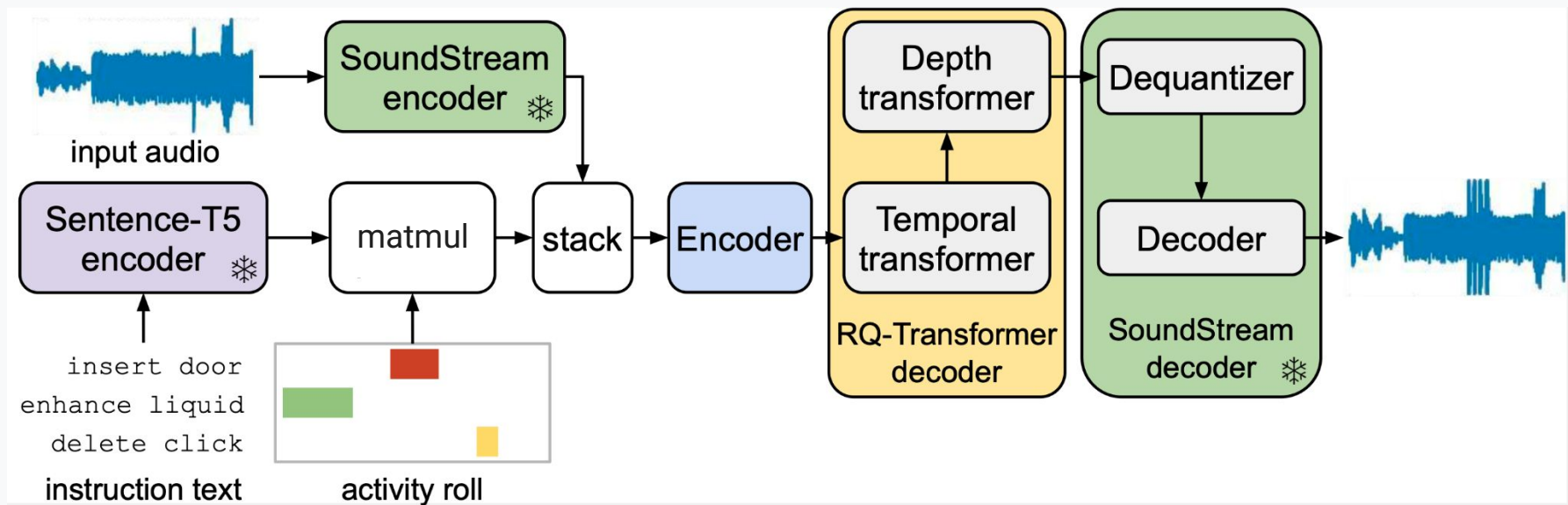


- Represents sound mixture as a collection of **discrete events**
- Each row associated with free-text **label** (Action + Class)
- ***Could*** be generated by an Audio Event Classifier

Contribution: **Enhancement**

- Enhancement combines **sound separation** and **generation**
 - Identify low-level target event in the input
 - Regenerate missing details to synthesize output
- How to specify “Enhancement”?
 - Constant **gain** applied to input?
 - Constant output target **level** (e.g. 10 dB) regardless of input level
 - (allow user to specify)

Recomposer Model

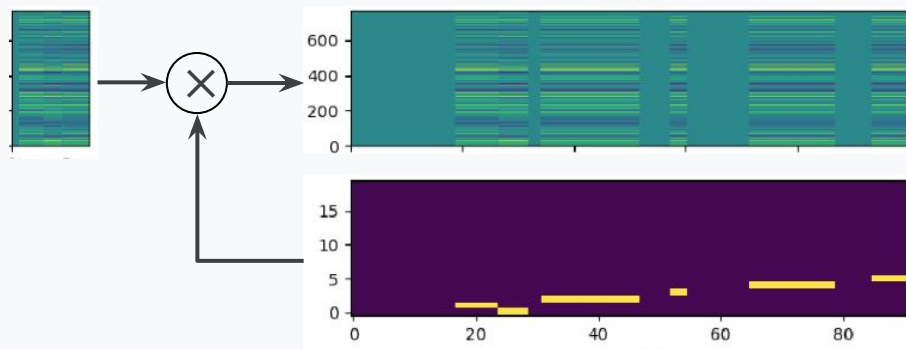
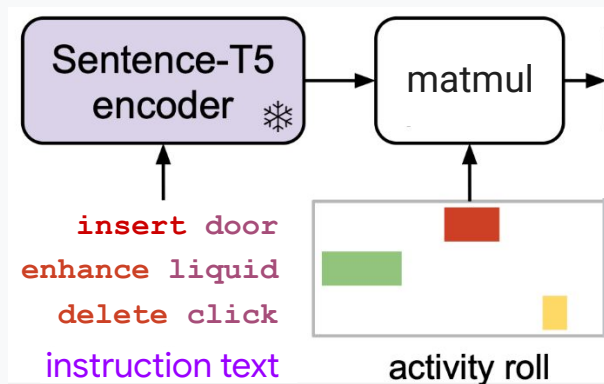


- Input sound **embedding + conditioning** encodes to conditioning **embedding**
- Autoregressive encoder-decoder transformer trained to generate **RVQ tokens**

[Autoregressive Image Generation using Residual Quantization](#), Lee et al, CVPR 2022

Recomposer Model: Edit conditioning

- **Instruction text** is embedded as **general text**
 - **Action** and **Target class** confounded...
 - Permits semantic **adjacency**
- `matmul((embedding, event), (event, time))`
 $\rightarrow (embedding, time)$

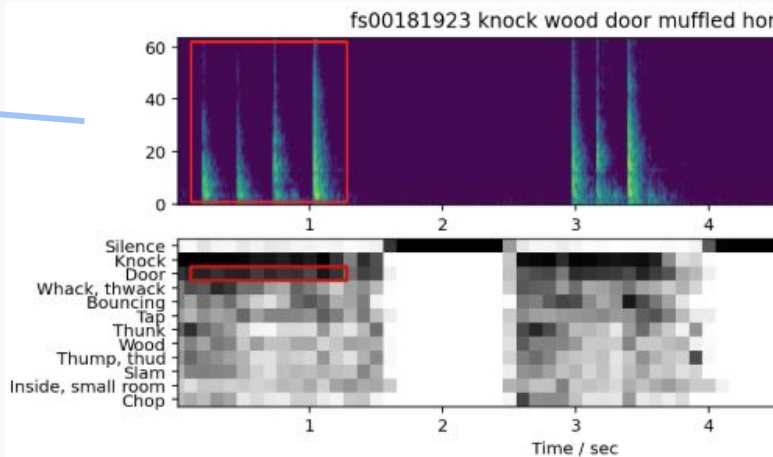
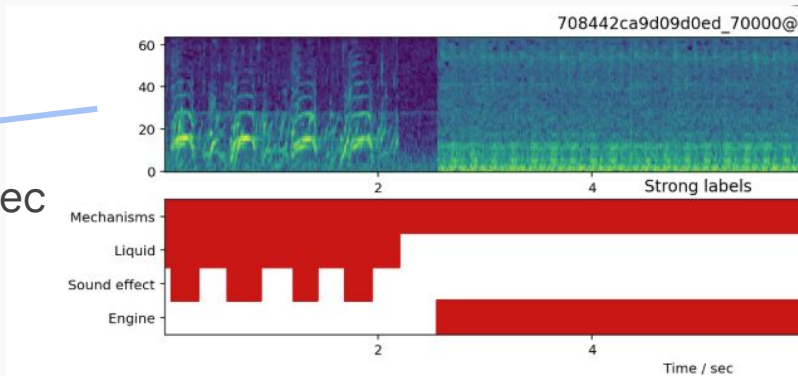


- All edits summed into
a **single edit-conditioning vector** per time step

Synthetic Target Data

Training examples as Background + Target

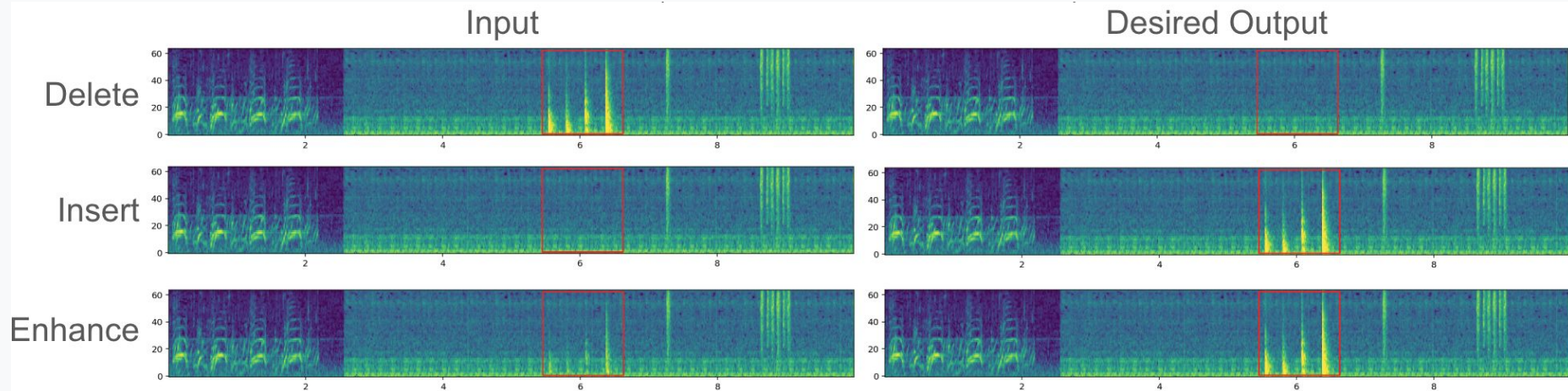
- **Background** from **AudioSet** (Strong Labeled) (10 sec clips)
 - Selected for Complexity (> 2 classes) and Time coverage (> 9 sec nonsilent)
 - 168k clips
- **Target** events from **Freesound**
 - Matched by tags + AudioSet classifier
 - Automatic trimming to isolated events of 0.2 .. 2 sec
 - 16k training events from 40 “event-like” AudioSet classes



Synthetic Target Data (cont'd)

- Synthetic data $\{input, output\}$ pairs formed by mixing Target + Background
- Fix **Target-to-Background Ratio (TBR)** vs. the overlapping background

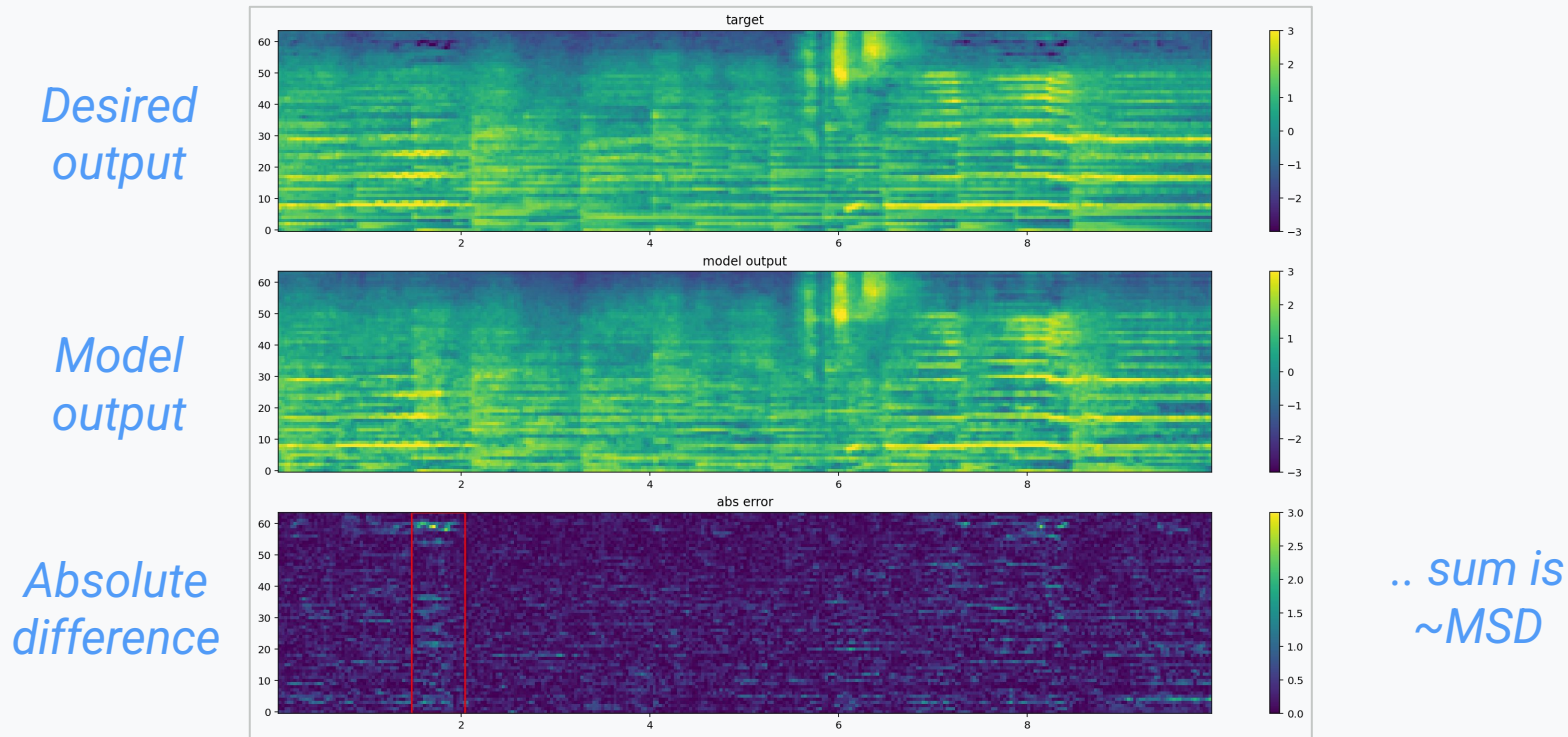
Three kinds of Action



Synthetic Target Data details

- Training: Generate **on-the-fly** *{input, output}* pairs
 - Two “Targets” per example
 - Each chosen from Enhance, Delete, Insert, or None (EDIN)
 - No time overlap between targets
 - Targets inserted at 10 (± 3) dB TBR
 - Loud enough to measure, still plausible
 - Enhancement inputs at -6 (± 3) dB TBR

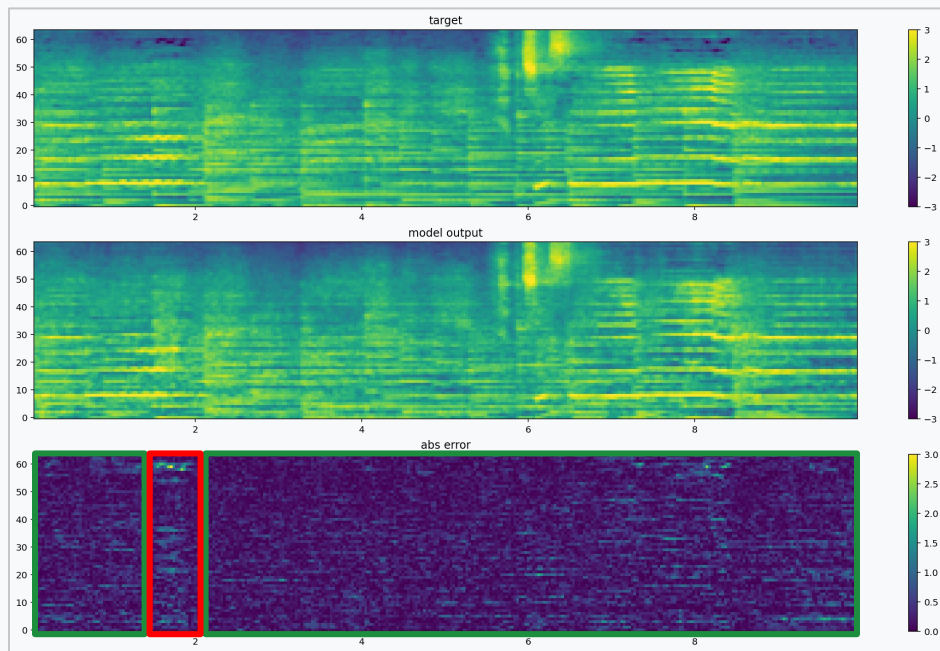
Metrics: Multi-Scale Spectral Distortion (MSD)



- Constant baseline “**model distortion**” – but differences are not audible

Metrics: Modified vs. Unmodified regions

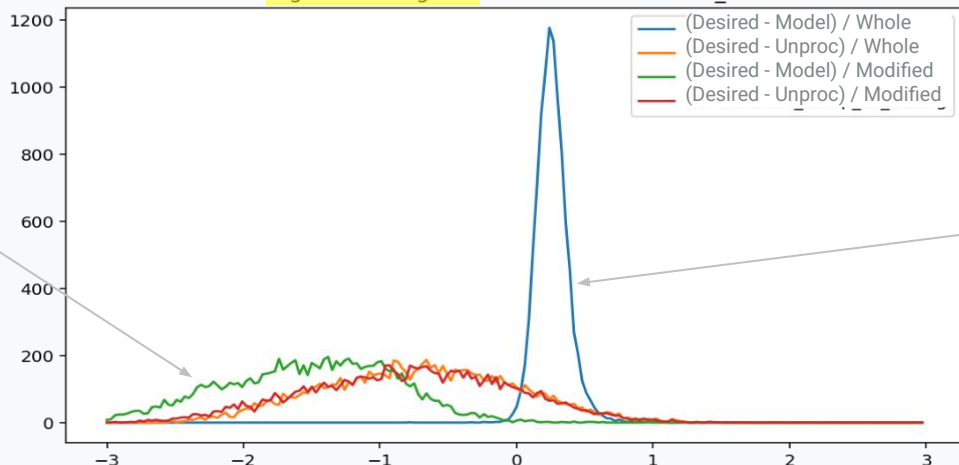
- Target event occupies ~10% of time
- “Model distortion” over remainder can dominate metrics
- **Framewise** metrics allow decomposition by time range
 - Use ground-truth target region to calculate separate metrics for **modified** & **unmodified** regions



Metrics: Modified vs. Unmodified regions

- Look at **histograms** of metrics across eval set
 - Compare $(\text{Desired} - \text{Model_output})$ and $(\text{Desired} - \text{Unprocessed})$
 - Calculate metrics over `Whole_clip` vs. `Modified_region`
- “Model distortion” overwhelms $(\text{Desired} - \text{Model_output}) / \text{Whole_clip}$
Restricting to `Modified_region` reveals benefit:

log MSD histogram: Delete xid112478261_diff



*Modified_region metrics:
Model_output **better**
than Unprocessed*

*Whole_clip metrics:
Model_output **worse**
than Unprocessed*

Metrics: Classifier KL Divergence (KLD)

Results: Varying Input Level for Enhancement

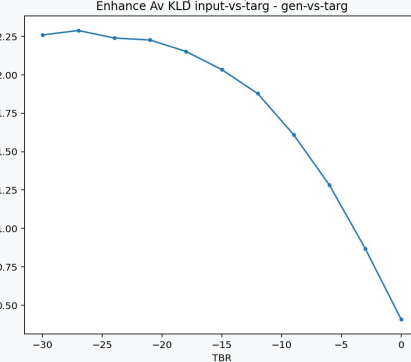
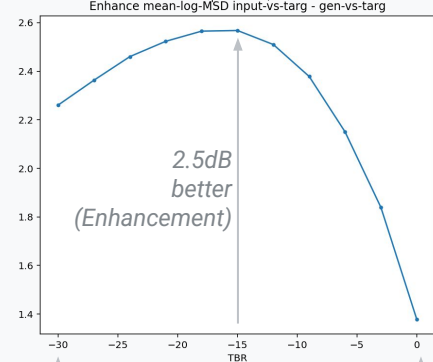
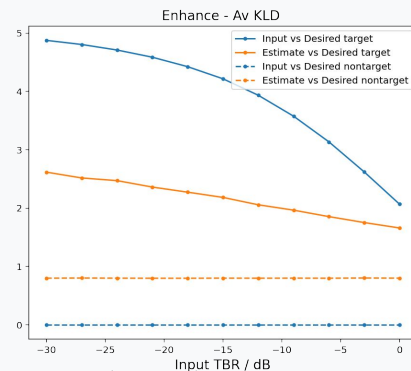
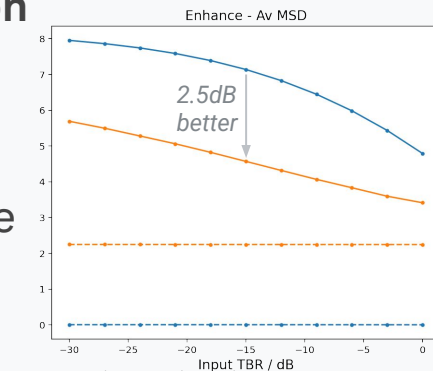
- Focus on generation-to-separation **transition** in Enhancement

- Training:

- One Enhancement target per example
- Inputs ranged over -30 to 0 dB TBR
- Output always at +15 dB (minimize uncertainty)

- Evaluation:

- Metrics improve with TBR
- Improvement-over-unprocessed peaks at ~ -15 dB TBR (for MSD)
- (KLD has no peak)



mostly
generation

mostly
separation

Results: Overall Performance

Region	Signal	Delete		Insert		Enhance	
		MSD	KLD	MSD	KLD	MSD	KLD
Target	input	4.8	1.6	4.8	2.8	3.4	1.6
	estimate	2.5	0.5	5.1	1.9	2.6	0.9
Nontarget	input	0.0	0.0	0.0	0.0	0.0	0.0
	estimate	1.3	0.3	1.3	0.3	1.3	0.3

(lower is better)

- Nontarget **input** is perfect - distances of 0
- Nontarget **estimate** is “**model distortion**” - limitation of copying
- Target **estimate** minus **input** gives **improvement** from processing
 - For MSD, **Delete** does well, **Insert** is made worse (specific output is unknown)
 - For KLD, **Insert** reveals benefit (because target class is specified)
 - **Enhance** results are for ~ -6 dB TBR inputs, limited headroom

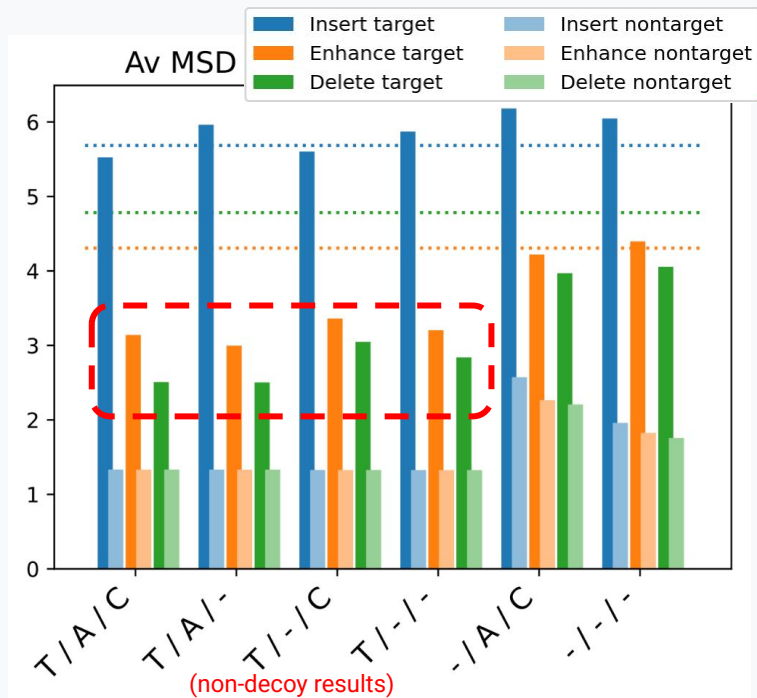
Conditioning Ablation

- Results show improvement - we're done?
 - Or is there more insight?
- Conditioning has 3 components
 - **Timing** (Event Roll)
 - **Action** (Delete/Insert/Enhance)
 - **Class** (Description of target)
- Train 6 models with partial conditioning
- Evaluate with **Decoy** examples
 - Minimize cues in input



Conditioning Ablation

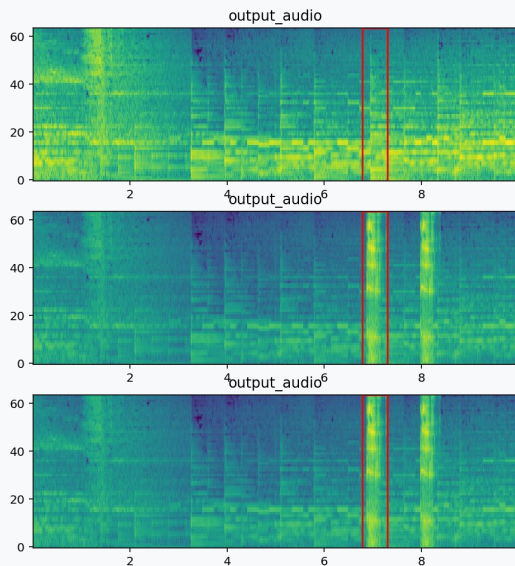
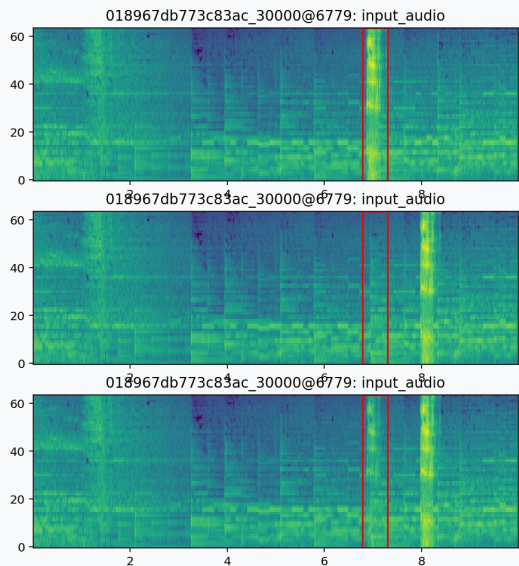
- Initial results: **Action** (and **Class**) have little impact
 - Model can infer them from input?



Decoy Data

Each sample input has **one event** (or decoy) at 10 dB TBR

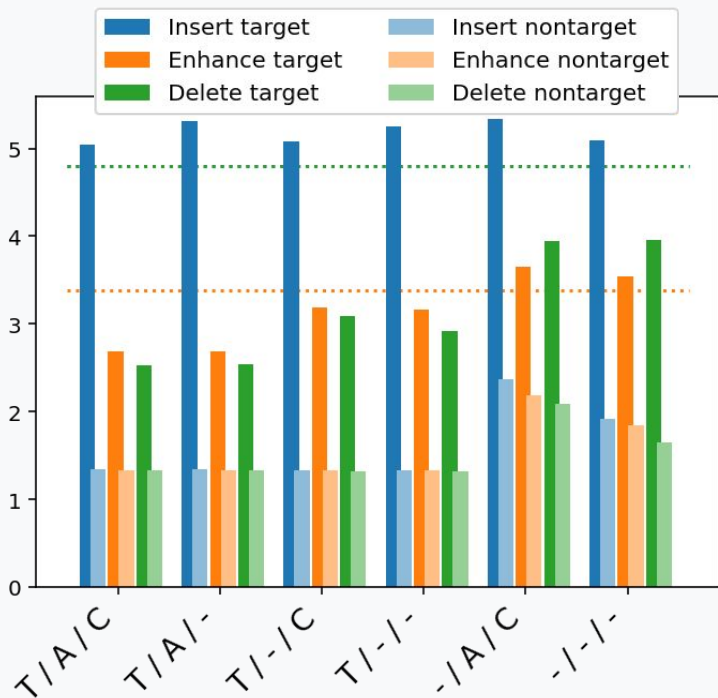
- **Delete:** Output has no event
- **Insert:** Output has a second event (event in input is **decoy**)
- **Enhance:** Input has decoy plus -6dB TBR target event



⇒ Model cannot guess action from input

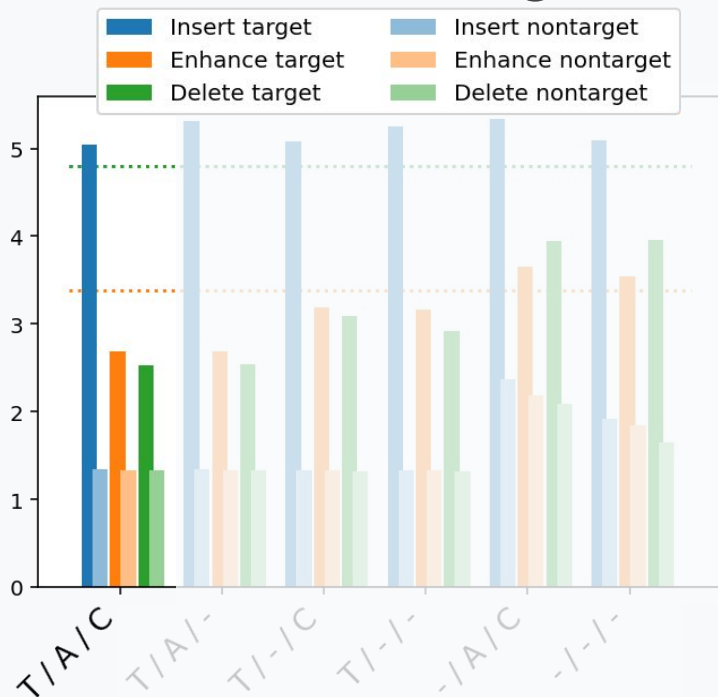
- Model was trained with 0 .. 2 edits per sample

Results: MSD



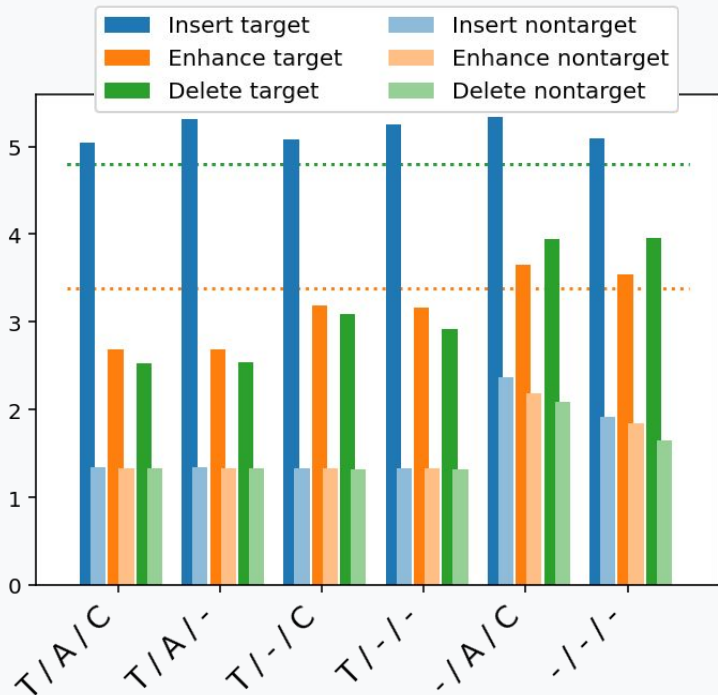
← Timing
Action
Class

Results: MSD - Full-Conditioning



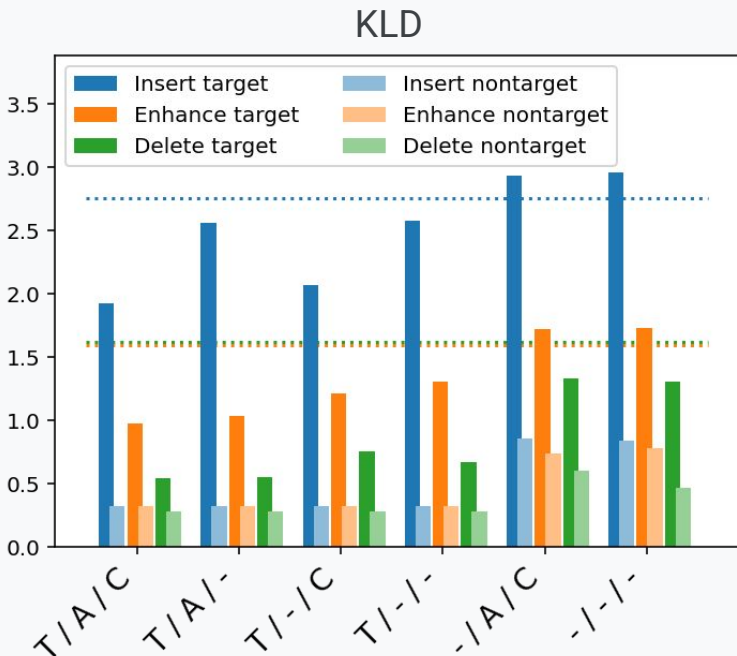
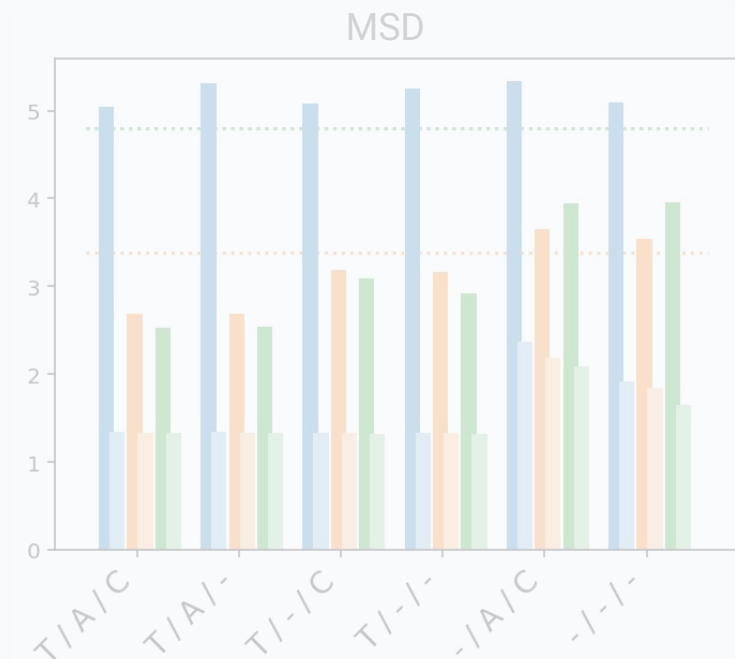
- Unprocessed baseline (dotted): **Insert & Delete** same (4.8), **Enhance** better (3.3)
- Unmodified region (nontarget, pale bars): “Model distortion” floor (1.3)
- Modified region (target, dark bars): **Delete ~ Enhance** (2.5), **Insert** much worse (5.0)

Results: MSD - Conditioning Ablation



- Remove Class (T / A / -) → **Insert** worse, **Delete & Enhance** unchanged?
- Remove Action (T / - / C) → **Delete & Enhance** worse (confused?)
- Remove Timing (- / A / C) → **Delete, Enhance** (and nontarget) worse

Results: KLD



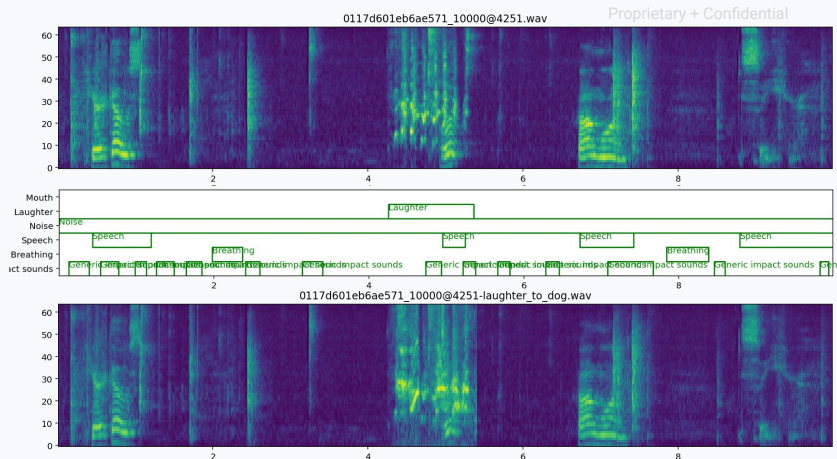
- Unprocessed baseline (dotted): **Delete** \approx **Enhance**, but **Insert** much worse
- Processed: **Delete** better than **Enhance**
- **Insert** benefits from Class









Future Work

- Richer **control** of generation
 - More structured attributes (loudness, pitch, reverberation)
 - Richer text-to-audio generation
 - training data?
- Broader conditioning e.g. **Video**
 - Audio-Video joint generation
- Sound **Transformation** ...

Sound Transformation via **Cross-Class Enhancement**

- Input audio contains a given `source_class`
- Run an enhancement-only `Recomposer` model to `'enhance <destination_class>'` at times when the **source_class** occurs.
 - The model was never trained to do this



Clip ID	0239dc6ce0480dc9_30000@4598	0253eeff2c4b4f68_230000@5512	0117d601eb6ae571_10000@4251	33ad41049a6fbf1f_30000@3207
Source class	Human locomotion	Cough	Laughter	Ring
Destination class	Digestive	Dog	Dog	Bird
Input audio				
Output audio				 Google

Conclusions

- Autoregressive Encoder-Decoder models can **edit scene details**
 - **Event roll** as a precise way to specify timing
- **Ablations** reveal complex interactions
 - Each part of the conditioning has a different effect
- Future work:
 - Richer **control** of generated sound events
 - Additional conditioning, e.g. **video**