

A. RELATED WORKS

A.1. Composition of Audio Processors

Most audio processors are designed to modify some specific properties of their input signals, e.g., magnitude response, loudness dynamics, and stereo width. As such, combining multiple processors is a common practice to achieve the full effect. Following the main text, we will use the terminology “graph” to represent this composition, although some previous works considered simple structures that allow a more compact form, e.g., a sequence. Now, we outline the previous attempts that tried to estimate the processing graph G and/or its parameters \mathbf{P} from reference audio. For example, if the references are one or more dry sources and a wet mixture, this task becomes reverse engineering [1, 2, 3, 4]. In terms of the prediction targets, some works [5, 6, 4, 7, 8] fixed the graph and estimated only the parameters. Others [1, 2, 3, 9, 10] tried also to predict the graph. These works differ in task and domains, processor sets, search space of graphs, and estimation methods. For example, Lee et al. [1] tackled singing voice effects and drum mixing estimation from wet and optional dry source(s). With 33 processor types where some of them were multiple-input multiple-output (MIMO) systems, they generated synthetic directed acyclic graphs (DAGs) and trained a transformer-based graph and parameter predictor in a supervised manner. These are clearly different from our work. We tackled the music mixing domain, which has much larger graph sizes but with a restricted graph structure based on the mixing consoles (more details on the resulting search spaces are in Appendix A.3). We used fewer processor types, but they were all single-input single-output (SISO) and differentiable. We ran joint pruning and parameter optimization of the hybrid mixing consoles via gradient descent (without any neural network). Such differences between the works are summarized and highlighted in Table 1. For a unified viewpoint, we use a probabilistic framework to denote the tasks, e.g., $p(\mathbf{P}|\mathbf{S})$ for automatic mixing, albeit many works do not involve any generative modeling or sampling.

A.2. Differentiable Signal Processing

Differentiable processor — Exact implementation or approximation of processors in an automatic differentiation framework, e.g., pytorch [11], enables parameter optimization via gradient descent. Numerous efforts have been focused on converting existing audio processors to their differentiable versions [4, 12, 13, 14, 15, 16]; refer to the recent review [17] and references therein for more details. In many cases, these processors are combined with neural networks, whose computation is done in GPU. Thus, converting the audio processors to be “GPU-friendly” has been an active research topic. For example, for a linear time-invariant (LTI) system with a recurrent structure, we can sample its frequency response to approximate its infinite impulse response (IIR) instead of directly running the recurrence; the former is faster than the latter [15, 16]. However, it is nontrivial to apply a similar trick to nonlinear, recurrent, or time-varying processors. Typically, further simplifications and approximations are employed, e.g., replacing the nonlinear recurrent part with an IIR filter [12] or assuming frame-wise LTI to a sample-level linear time-varying system [14]. Sometimes, we can only access input and output measurements. Then, one can approximate the gradients with finite difference methods [8, 18] or pretraining an auxiliary neural network that mimics the target processors [5]. In the literature, these tricks are also referred to as “differentiable;” hence, it is rather an umbrella term encompass-

ing all possible methods that obtain the output signals or gradients within a reasonable time. Nevertheless, we limit our focus to the implementations in the automatic differentiation framework.

Audio processing graph — Now, consider a composition of multiple differentiable processors; the entire graph remains differentiable due to the chain rule. However, the following practical considerations remain. If we fix the processing graph prior to the optimization and the graph size is relatively small, we can implement the “differentiable graph” following the existing implementations [7, 6]. That is, we compute every processor one by one in a pre-defined topological order. However, we have the following additional requirements. First, the pruning changes the graph during the optimization. Therefore, our implementation must take a graph and its parameters along with the source signals as input arguments for every forward pass. Note that this feature is also necessary when training a neural network that predicts the parameters of any given graph [1]. Second, the size of our graphs is much larger than the ones from previous works [7, 6, 9, 1]. In this case, the one-by-one computation severely bottlenecks the computation speed. Therefore, we derived a flexible and efficient graph computation algorithm (i) that can take different graphs for each forward pass as input and (ii) performs batched processing of multiple processors within a graph, utilizing the parallelism of GPUs. Finally, we note that other than the differentiation with respect to the input signals \mathbf{S} and parameters \mathbf{P} , one might be interested in differentiation with respect to the graph structure G . The proposed pruning method performs this to a limited extent; deletion of a node v_i is a binary operation that modifies the graph structure. We relaxed this to a continuous dry/wet weight w_i and optimized it with the audio loss L_a and regularization L_p .

A.3. Graph Search

Several independent research efforts in various domains exist that aim to find graphs that satisfy certain requirements. First, neural architecture search (NAS) aims to find a neural network architecture that achieves an improved performance. In this case, the search space consists of graphs where the nodes (or edges according to the chosen notation) represent the primitive neural network layers. One very close work to ours is differentiable architecture search (DARTS) [19], which relaxed the choice of layer type to a categorical distribution and optimized it via gradient descent. Theoretically, our method can be naturally extended to follow this approach; we only need to change our 2-way choice (prune or not) to $(K+1)$ -way (bypass or select one of K possible processors). While this search is more flexible and general, this greatly increases the computational cost, as we must compute all K processor types for every node. Second, the generation/design of a molecule is another popular topic. One dominant approach is to use reinforcement learning (RL) to generate molecule graphs that achieve the desired properties [20]. However, applying this method to our task has a risk of obtaining nontrivial mixing graphs that is hard to interpret for practitioners.

B. DIFFERENTIABLE AUDIO PROCESSORS

For each processor, $u[n]$ and $y[n]$ denote its stereo input and output, respectively. We write the discrete-time index with n . p with a subscript denotes a real-valued parameter vector, which could be a concatenation of multiple parameters. Each subscript x denotes one of the following channels: left l, right r, mid m, and side s.

[6]	<i>Task & domain</i>	Sound matching $p(\mathbf{P} x)$. The synthesizer parameters \mathbf{P} were estimated to match the reference (target) audio x .
	<i>Processors</i>	Oscillators, envelope generators, and filters that allow parameter modulation as an optional input.
	<i>Graph</i>	Any pre-defined DAG. For example, a subtractive synthesizer (2 oscillators, 1 amplitude envelope, and 1 lowpass filter were used for the evaluation).
	<i>Optimization</i>	Trained a single neural backbone for the reference encoding, followed by multiple prediction heads for the parameters. Optimized with a parameter loss and spectral loss, where the latter is calculated with every intermediate output.
[10]	<i>Task & domain</i>	Sound matching $p(\mathbf{P} x)$. An FM synthesizer matches recordings of monophonic instruments (violin, flute, and trumpet). Estimates parameters of an operator graph G that is empirically searched & selected.
	<i>Processors</i>	Differentiable sinusoidal oscillators, each used as a carrier or modulator, pre-defined frequencies. An additional FIR reverb is added to the FM graph output for post-processing.
	<i>Graph</i>	Directed and acyclic (no feedback FM), at most 6 operators. Different graphs for different target instruments.
	<i>Optimization</i>	Trained a convolutional neural network that estimates envelopes from the target loudness and pitch.
[9]	<i>Task & domain</i>	Sound matching $p(G, \mathbf{P} x)$. Similar setup to the above [10] but with an additional estimation of the operator graph G .
	<i>Processors</i>	Identical to [10], except for the frequency ratio that can be searched.
	<i>Graph</i>	A subgraph of a supergraph, which resembles a multi-layer perceptron (modulator layers followed by a carrier layer).
	<i>Optimization</i>	Trained a parameter estimator for the supergraph and found the appropriate subgraph G with an evolutionary search.
[2]	<i>Task & domain</i>	Reverse engineering $p(G, \mathbf{P} s, y)$ of an audio effect chain from a subtractive synthesizer (commercial plugin).
	<i>Processors</i>	5 audio effects: compressor, distortion, equalizer, phaser, and reverb. Non-differentiable implementations.
	<i>Graph</i>	Chain of audio effects generated with no duplicate types (therefore 32 possible combinations) and random order.
	<i>Optimization</i>	Trained a next effect predictor and parameter estimator in a supervised (teacher-forcing) manner.
[3]	<i>Task & domain</i>	Blind estimation $p(G y)$ and reverse engineering $p(G s, y)$ of guitar effect chains.
	<i>Processors</i>	13 guitar effects, including non-linear processors, modulation effects, ambience effects, and equalizer filters.
	<i>Graph</i>	A chain of guitar effects. Maximum 5 processors and a total of 221 possible combinations.
	<i>Optimization</i>	Trained a convolutional neural network with synthetic data to predict the correct combination.
[5]	<i>Task & domain</i>	Automatic mixing $p(\mathbf{P} \mathbf{S})$. Estimated parameters of fixed processing chains from source tracks ($K \leq 16$).
	<i>Processors</i>	7 differentiable SISO processors, where 4 (gain, polarity, fader, and panning) were implemented exactly. The remaining 3 (equalizer, compressor, and reverb) were approximated with a single pretrained neural network.
	<i>Graph</i>	Tree structure: applied a fixed chain of the 7 processors for each track, and then summed the chain outputs altogether.
	<i>Optimization</i>	Trained a parameter estimator (convolutional neural network) with a spectrogram loss end-to-end.
[8]	<i>Task & domain</i>	Reverse engineering of music mastering $p(\mathbf{P} s, y)$
	<i>Processors</i>	A multi-band compressor, graphic equalizer, and limiter. Gradient approximated with a finite difference method.
	<i>Graph</i>	A serial chain of the processors.
	<i>Optimization</i>	Optimized parameters with gradient descent.
[1]	<i>Task & domain</i>	Blind estimation $p(G, \mathbf{P} y)$ and reverse engineering $p(G, \mathbf{P} \mathbf{S}, y)$. Estimates both graph and its parameters, for singing voice effect ($K = 1$) or drum mixing ($K \leq 6$).
	<i>Processors</i>	A total of 33 processors, including linear filters, nonlinear filters, and control signal generators. Some processors are MIMO systems, e.g., allowing auxiliary modulations. Non-differentiable implementations.
	<i>Graph</i>	Complex DAG; splits (e.g., multi-band processing) and merges (e.g., sum and modulation). 30 processors max.
	<i>Optimization</i>	Trained a convolutional neural network-based reference encoder and a transformer variant for graph decoding and parameter estimation. Both were jointly trained via direct supervision of synthetic graphs (e.g., parameter loss).
[4]	<i>Task & domain</i>	Reverse engineering $p(\mathbf{P} \mathbf{S}, y)$ of music mixing. Estimated parameters of a fixed chain for each track.
	<i>Processors</i>	6 differentiable SISO processors: gain, equalizer, compressor, distortion, panning, and reverb.
	<i>Graph</i>	A chain of 5 processors (all above types except the reverb) for each dry track (any other DAG can also be used). The reverb is used for the mixed sum.
	<i>Optimization</i>	Parameters were optimized with spectrogram loss end-to-end via gradient descent.
Ours	<i>Task & domain</i>	Reverse engineering $p(G, \mathbf{P} \mathbf{S}, y)$ of music mixing. Estimated a chain of processors and their parameters for each track and submix where $K \leq 130$.
	<i>Processors</i>	7 differentiable SISO processors: gain/panning, stereo imager, equalizer, reverb, compressor, noisegate, and delay.
	<i>Graph</i>	A tree of processing chains with a subgrouping structure (any other DAG can also be used). Processors can be omitted but should follow the fixed order.
	<i>Optimization</i>	Joint estimation of the soft masks (dry/wet weights) and processor parameters. Optimized with the spectrogram loss (and additional regularizations) end-to-end via gradient descent. Accompanied by hard pruning stages.

Table 1: A brief summary and comparison of previous works on estimation of compositional audio signal processing.

B.1. Gain/panning

We use simple channel-wise constant multiplication. Its parameter vector $p_g \in \mathbb{R}^2$ is in log scale, so we apply exponentiation before multiplying it to the stereo signal.

$$y[n] = \exp(p_g) \cdot u[n]. \quad (1)$$

B.2. Stereo Imager

We multiply the side signal (left minus right) with a gain parameter $p_s \in \mathbb{R}$ to control the stereo width. The mid and side outputs are

$$y_m[n] = u_l[n] + u_r[n], \quad (2a)$$

$$y_s[n] = \exp(p_s) \cdot (u_l[n] - u_r[n]). \quad (2b)$$

Then, we convert the mid/side output to a stereo signal as follows, $y_l[n] = (y_m[n] + y_s[n])/2$ and $y_r[n] = (y_m[n] - y_s[n])/2$.

B.3. Equalizer

We employ a zero-phase FIR filter. Considering its log-magnitude as a parameter p_e , we compute inverse FFT (IFFT) of the magnitude response and multiply it with a Hann window $v^{\text{Hann}}[n]$. As a result, the length- N FIR is given as

$$h_e[n] = v^{\text{Hann}}[n] \cdot \frac{1}{N} \sum_{k=0}^{N-1} \exp p_e[k] \cdot w_N^{kn} \quad (3)$$

where $-(N+1)/2 \leq n \leq (N+1)/2$ and $w_N = \exp(j \cdot 2\pi/N)$. We compute the final output by applying the same FIR to both the left and right channels as follows,

$$y_x[n] = u_x[n] * h_e[n] \quad (x \in \{l, r\}). \quad (4)$$

We set the FIR length to $N = 2047$. Therefore, the parameter p_e has a size of 1024. Note that we could use a parametric equalizer [15] as an alternative for more compact parameters. However, we used this FIR filter due to its simplicity and fast computation.

B.4. Reverb

We use a variant of filtered noise models [21, 7]. First, we create 2 seconds of uniform noises, $u_m[n]$ and $u_s[n]$. Next, we multiply each noise's STFT $U_x[k, m]$ with a magnitude mask $M_x[k, m]$.

$$H_x[k, m] = U_x[k, m] \odot M_x[k, m] \quad (x \in \{m, s\}). \quad (5)$$

Here, k and m denote frequency and time frame index. Each mask is parameterized with an initial log-magnitude $H_x^0[k]$ and an absorption filter $H_x^\Delta[k]$ as follows,

$$M_x[k, m] = \exp(H_x^0[k] + (m-1)H_x^\Delta[k]). \quad (6)$$

Next, we convert the masked STFTs to the time-domain responses, $h_m[n]$ and $h_s[n]$. We obtain the desired FIR $h_r[n]$ by converting the mid/side to stereo. We apply channel-wise convolutions to the input $u[n]$ and get the output signal $y[n]$. The FFT and hop lengths are set to 384 and 192, respectively. The parameter vector p_r has the size of 768 (2 channels, each with 2 filters, and 192 magnitude values for each filter). We briefly discuss some alternatives to this processor. First, we could use other artificial reverberation models [16]. But again, we chose this STFT-based FIR due to its simplicity and fast computation. Or, we could use a single FIR to represent both reverb and multitap delay since they are linear time-invariant [22]. However, we decided to separate those two for interpretable, controllable, and more compact representations.

B.5. Compressor

We implement the canonical feed-forward digital compressor [23]. First, for a given input audio, we sum the left and right channels to obtain a mid signal $u_m[n]$. Then, we calculate its energy envelope $G_u[n] = \log g_u[n]$ where

$$g_u[n] = \alpha[n]g_u[n-1] + (1-\alpha[n])u_m^2[n]. \quad (7)$$

Here, the coefficient $\alpha[n]$ is typically set to a different constant for an “attack” (increasing) and “release” (decreasing) phase:

$$\alpha[n] = \begin{cases} \alpha^{\text{att}} & g_u[n] > g_u[n-1], \\ \alpha^{\text{rel}} & g_u[n] \leq g_u[n-1]. \end{cases} \quad (8)$$

As this part (also known as ballistics) bottlenecks the computation speed in GPU, following the recent work [12], we restrict the constants to the same value: $\alpha = \alpha^{\text{att}} = \alpha^{\text{rel}}$. By doing so, Equation 7 simplifies to a one-pole IIR filter, whose length- N FIR approximation can be computed in parallel using the frequency-sampling method [15] and applying IFFT to the sampled response:

$$h^{\text{env}}[n] \approx \frac{1}{N} \sum_{k=0}^{N-1} \frac{1-\alpha}{1-\alpha w_N^{-k}} w_N^{kn}. \quad (9)$$

Therefore, the energy envelope can be simply computed as

$$g_u[n] \approx h^{\text{env}}[n] * u^2[n]. \quad (10)$$

Next, we compute the compressed energy envelope $G_y[n]$. We use a quadratic knee, interpolating the compression and the bypass region. For a given threshold T and half of the knee width W ,

$$G_y[n] = \begin{cases} G_y^{\text{above}}[n] & G_u[n] \geq T + W, \\ G_y^{\text{mid}}[n] & T - W \leq G_u[n] < T + W, \\ G_y^{\text{below}}[n] & G_u[n] < T - W \end{cases} \quad (11)$$

where, for a given compression ratio R , each term is

$$G_y^{\text{above}}[n] = T + \frac{G_u[n] - T}{R}, \quad (12a)$$

$$G_y^{\text{mid}}[n] = G_u[n] + \left(\frac{1}{R} - 1\right) \frac{(G_u[n] - T + W)^2}{4W}, \quad (12b)$$

and $G_y^{\text{below}}[n] = G_u[n]$. Finally, we can compute the output as

$$y_x[n] = \exp(G_y[n] - G_u[n]) \cdot u_x[n] \quad (x \in \{l, r\}). \quad (13)$$

The scalar parameters introduced above, α , T , W , and R , are concatenated and used as a parameter vector $p_c \in \mathbb{R}^4$. Our implementation is almost identical to the recently introduced differentiable compressor [12]. However, we smoothed to the energy $u_m^2[n]$ instead of the gain reduction $G_y[n] - G_u[n]$ as the former showed a slightly better matching performance in the initial experiments. Since the approximated ballistics (Equation 9) could be one cause of the reduced modeling capability of loudness dynamics, it might be desirable to try alternatives that allow the different attack and release coefficients [18, 24].

B.6. Noisegate

Its implementation is the same as the compressor above, except for the gain computation part. We set $G_y^{\text{above}}[n] = G_u[n]$ and

$$G_y^{\text{mid}}[n] = G_u[n] + (1-R) \frac{(G_u[n] - T - W)^2}{4W}, \quad (14a)$$

$$G_y^{\text{below}}[n] = T + R(G_u[n] - T). \quad (14b)$$

B.7. Multitap Delay

We consider a 2 seconds of delay effect with at most one delay d_m at every 100ms (hence the number of delay taps is $M = 20$). Each delay is filtered with a 39-tap zero-phase FIR $c_m[n]$ parameterized in the same way as the equalizer. We use independent delays and their filters for the left and right channels, but we ignore this in the following equations for simplicity. Under this setting, the multitap delay's single-channel FIR $h_{\text{d}}[n]$ is given as follows,

$$h_{\text{d}}[n] = \sum_{m=1}^M c_m[n] * \delta[n - d_m] \quad (15)$$

where $\delta[n]$ is an unit impulse. Here, we aim to optimize each delay length $d_m \in \mathbb{N}^+$, which is discrete, using gradient descent. To this end, we exploit the fact that each delay $\delta[n - d_m]$ corresponds to a complex sinusoid in the frequency domain. Recent work showed that the sinusoid's angular frequency z_m can be optimized with the gradient descent when we allow it to be inside of the unit disk [25]. In short, for each complex parameter $z_m \in \mathbb{C}$ where $|z_m| \leq 1$, we compute a surrogate damped sinusoid and apply IFFT to it.

$$\delta[n - d_m] \approx \frac{1}{N} \sum_{k=0}^{N-1} z_m^k w_N^{kn}. \quad (16)$$

This relaxed FIR is not exactly the discrete delay signal; precisely, it is a lowpassed version of an aliased sinc kernel. Hence, we use it only for backpropagation with the straight-through technique [26]. Also, we normalize each gradient and regularize the parameter z_m to be closer to the unit circle. We empirically observed that this improves the performance. To summarize, each complex conjugate gradient is modified as follows and used for the optimization:

$$\frac{\partial L}{\partial z_m^*} \leftarrow \text{sgn}\left(\sum_n \frac{\partial L}{\partial h_{\text{d}}[n]} \frac{\partial \tilde{h}_{\text{d}}[n]}{\partial z_m^*}\right) + \gamma(|z_m| - 1) \text{sgn}(z_m^*) \quad (17)$$

where $\tilde{h}_{\text{d}}[n]$ is the surrogate FIR obtained with the continuous delay, $\text{sgn}(z) = z/|z|$, and γ is a regularization strength set to 0.01. This multitap delay has parameter p_{d} of size 880 (20 delays each for the left and right channel; each delay with a real and imaginary part of the sinusoid's angular frequency and 20 log-magnitude bins of the FIR filter). Note that we restricted the delay lengths to be in separate time intervals; we failed to optimize this processor without this constraint. Resolving this is left as a future work.

C. AUDIO PROCESSING GRAPHS IN GPU

We slightly deviate from the main text and consider a more general setup. We do not distinguish between the processors and auxiliary modules; both are referred to as nodes. We omit all the weights w for simplicity. Then, each node's output y_i is given as follows,

$$y_i = f_i(u_i, p_i), \quad u_i = \sum_{j \in \mathcal{N}^+(i)} y_j \quad (18)$$

where $\mathcal{N}^+(i)$ denotes a set of source node indices for the node v_i . In short, each output y_i is obtained by gathering the input signals, aggregating those, and processing the sum with the parameters (the latter two can be omitted according to the node type). Furthermore, we allow disconnected graphs and view a batch of multiple graphs as a single large graph. Thus, it is enough to consider only a single graph case. The remaining structural restrictions to graphs are (i) being acyclic and (ii) having only SISO systems.

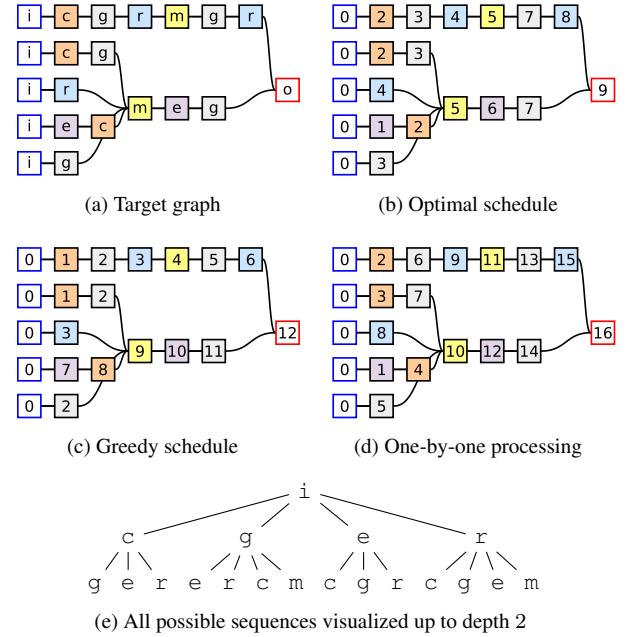


Figure 1: Various batched processing schedules. For each schedule 1b-1d, the processing orders are shown inside the nodes.

C.1. Batched Processing

We may obtain the final graph output \mathbf{Y} by simply processing each node one by one. Nevertheless, we can accelerate the computation via parallelism. Consider a subset sequence $V_0, \dots, V_N \subset V$ that satisfies the followings,

- (i) It is a *partition*: $\cup_n V_n = V$ and $V_n \cap V_m = \emptyset$ for $n \neq m$.
- (ii) The sequence must be *causal*: there is no path from $u \in V_n$ to $v \in V_m$ if $n \geq m$.
- (iii) Each subset V_n is *homogeneous*: it has only one type t_n .

Then, we can compute a batch of output signals \mathbf{Y}_n of each subset V_n sequentially, from $n = 0$ to N . Consequently, we reduce the number of the gather-aggregate-process iterations from $|V|$ to N . Figure 1 shows an example. For a graph with $|V| = 21$ nodes (1a), we can find a sequence with $N = 9$ (1b). Intuitively, this batched processing is effective for a graph with fewer types and a structure that resembles the (pruned version of) mixing console, e.g., having the identical processor order for each track.

C.2. Node Type Scheduling

To maximize the parallelism, we want to find the shortest sequence. This is a variant of the scheduling problem [27]. First, we always choose a maximal subset V_i when the type t_i is fixed. This makes the subset sequence equivalent to a type string, e.g., `i e c g r m e g r o` for 1b. We also choose the initial and the last subset, V_0 and V_N , to collect all input and output nodes, respectively. Since the search tree for the shortest sequence exponentially grows, the brute-force search is too expensive for most graphs (1e). Fortunately, we know the optimal schedule for the graphs in this paper as they follow the fixed ordering. If we do not have such restrictions or prior assumptions about the graph structure, we may try the greedy method that chooses a type with the largest number of computable nodes (1c). However, this usually results in a longer sequence.

Algorithm 1 Batch computation of audio processing graphs.

```

Input: Node types  $\mathbf{T}$ , edges  $\mathbf{E}$ , parameters  $\mathbf{P}$ , and inputs  $\mathbf{S}$ 
Output: Output signals  $\mathbf{Y}$  and (optional) intermediate signals  $\mathbf{U}$ 

1:  $\bar{\mathbf{T}}, N \leftarrow \text{ScheduleBatchedProcessing}(\mathbf{T}, \mathbf{E})$ 
2:  $\sigma \leftarrow \text{OptimizeNodeOrder}(\bar{\mathbf{T}}, \mathbf{T}, \mathbf{E})$ 
3:  $\mathbf{T}, \mathbf{E}, \mathbf{P} \leftarrow \text{Reorder}(\sigma, \mathbf{T}, \mathbf{E}, \mathbf{P})$ 
4:  $\mathbf{I}^G, \mathbf{I}^P, \mathbf{I}^A, \mathbf{I}^S \leftarrow \text{GetReadWriteIndex}(\bar{\mathbf{T}}, \mathbf{T}, \mathbf{E}, \mathbf{P})$ 
5:  $\mathbf{U} \leftarrow \text{Initialize}(\mathbf{S}, \mathbf{T})$ 
6: for  $n \leftarrow 1$  to  $N$  do
7:    $\bar{\mathbf{U}}_n \leftarrow \text{Gather}(\mathbf{U}, \mathbf{I}^G_n)$   $\triangleright \text{index\_select}$ 
8:    $\mathbf{U}_n \leftarrow \text{Aggregate}(\bar{\mathbf{U}}_n, \mathbf{I}^A_n)$   $\triangleright \text{scatter}$ 
9:    $\mathbf{P}_n \leftarrow \text{Gather}(\mathbf{P}_{[\bar{t}_n]}, \mathbf{I}^P_n)$   $\triangleright \text{slice}$ 
10:   $\mathbf{Y}_n \leftarrow \text{Process}(\bar{t}_n, \mathbf{U}_n, \mathbf{P}_n)$ 
11:   $\mathbf{U} \leftarrow \text{Store}(\mathbf{U}, \mathbf{Y}_n, \mathbf{I}^S_n)$   $\triangleright \text{slice}$ 
12: end for
13:  $\mathbf{Y} \leftarrow \mathbf{Y}_N$ 
14: return  $\mathbf{Y}, \mathbf{U}$ 

```

C.3. Implementation Details

Inputs — We represent a graph with a node type vector $\mathbf{T} \in \mathbb{N}^{|V|}$ and an edge index tensor $\mathbf{E} \in \mathbb{N}^{2 \times |E|}$. Its parameters are collected in a dictionary \mathbf{P} whose key is a type t and value is its corresponding parameter tensor given as $\mathbf{P}[t] \in \mathbb{R}^{|V_t| \times N_t}$ where $|V_t|$ and N_t are the number of nodes and parameters. All sources are stacked to a single tensor $\mathbf{S} \in \mathbb{R}^{K \times 2 \times L}$. We ensure that all the inputs, \mathbf{T} , \mathbf{E} , \mathbf{P} , and \mathbf{S} , share a node order. For example, a k^{th} source s_k must correspond to the first k^{th} input \downarrow in the type list \mathbf{T} . Likewise, an l^{th} type- t parameter $\mathbf{P}[t]_l \in \mathbb{R}^{N_t}$ must correspond to the first l^{th} type t in the type list \mathbf{T} .

Preprocessing — Algorithm 1 obtains the output $\mathbf{Y} \in \mathbb{R}^{|V_0| \times 2 \times L}$ from the prescribed inputs. First, we schedule the batched processings, resulting in a type list $\bar{\mathbf{T}} \in \mathbb{N}^N$ (1). Next, as the main batched processing loop (6-12) contains multiple memory read/writes, we calculate the node reordering σ that achieves contiguous memory accesses and improves the computation speed (2). This procedure allows memory accesses via `slice`, as shown in the comments of Algorithm 1. After permuting the graph tensors (3), we retrieve lists of tensor indices, \mathbf{I}^G , \mathbf{I}^P , \mathbf{I}^A , and \mathbf{I}^S , used for the read/writes in the main loop (4). Note that all these preprocessing steps (1-4) are done in CPU. Then, we create an intermediate output tensor $\mathbf{U} \in \mathbb{R}^{|V| \times 2 \times L}$ (5). As we put all the inputs to be the first partition V_0 , it can be initialized with simple concatenation: $\mathbf{U} = \mathbf{S} \oplus \mathbf{0}$.

Main loop — The remainings repeat batched processing and necessary read/writes (6-12). For each n^{th} iteration, we collect the previous outputs $\bar{\mathbf{U}}_n$ that are routed to the current partition nodes. We achieve this by accessing the intermediate tensor \mathbf{U} with the index \mathbf{I}^G_n with `index_select` (7). Then, we aggregate them using `scatter` operation if multiple edges are connected to some nodes (8). We can similarly obtain a parameter tensor \mathbf{P}_n with its corresponding index \mathbf{I}^P_n . Especially, our node reordering (3) makes this a simple `slice`, faster than the usual `index_select`. With the obtained input signals $\mathbf{U}_n \in \mathbb{R}^{|V_n| \times 2 \times L}$ and parameters $\mathbf{P}_n \in \mathbb{R}^{|V_n| \times N_t}$, we batch-compute the node outputs \mathbf{Y}_n (10). Then, we save them to the intermediate output tensor \mathbf{U} with the `slice` index \mathbf{I}^S_n (11) so that the remaining steps can access them as inputs. After the iteration, we have all node outputs saved in \mathbf{U} . The final graph outputs are given as $\mathbf{Y}_N \in \mathbb{R}^{|V_N| \times 2 \times L}$ since we set the last node partition V_N to collect all output nodes.

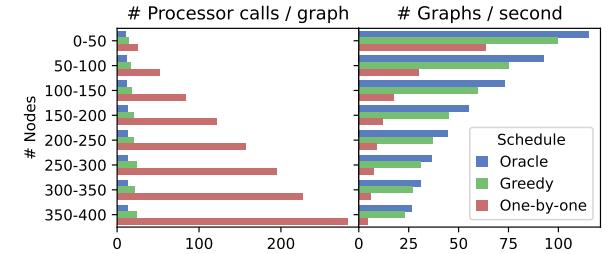


Figure 2: Benchmark results of different processing schedules.

C.4. Benchmark

We measured the efficiency of the schedules with the pruned graphs on a single RTX3090 GPU. Figure 2 reports the result.

- Our pruned graphs require at most 14 batched processings when scheduled correctly. The optimal schedule achieves 11.8 calls on average, while the greedy and one-by-one schedule report 17.5 and 77.6 calls, respectively. The latter especially increases the processor calls linearly to the graph size.
- As a result, the batched processing improves the speed across all graph sizes (especially for large ones). On average, the one-by-one computes 27.8 graph outputs per second. The greedy and optimal schedule achieve 67.2 and 81.2, respectively, indicating that improved scheduling methods for arbitrary graphs could be desirable for further research.
- Note that the speed-up from the batched processing is sublinear to the graph size. For graphs with 350 to 400 nodes, the optimal calls the processors 21.1 times less than the one-by-one, but the speedup is about 5.8 fold.

C.5. Extension

We can remove the SISO restriction and allow MIMO systems. A processor f_i with M inputs and N outputs processes the incoming signals as follows,

$$[y_{i1}, \dots, y_{iN}]^T = f_i([u_{i1}, \dots, u_{iM}]^T, p_i), \quad (19a)$$

$$u_{il} = \sum_{(j,k) \in \mathcal{N}^+(i,l)} y_{jk}. \quad (19b)$$

Here, $(j, k) \in \mathcal{N}^+(i, l)$ denotes a connection from node j channel k to node i channel l . In a graph viewpoint, this forms a multigraph that allows multiple edges between two nodes. The only remaining restriction (to the graphs) is being acyclic. Implementation-wise, we need to introduce an edge channel type tensor, replace the intermediate \mathbf{U} with a node-channel flattened tensor, and reshape the processor input/outputs, \mathbf{U}_n and \mathbf{Y}_n , appropriately to gather from and store to the intermediate \mathbf{U} .

D. DRY/WET PRUNING ALGORITHM

Algorithm 2 describes the details of the dry/wet method.

- For a simpler description, we modified the initialization part to include per-type node sets and weights, as in line 6-10.
- The termination condition is given in line 11.
- The trial candidate sampling is implemented in line 12-13.
- The candidate pool update is expanded to separately handle the trial successes and failures, shown in line 18 and 20-25.

	MedleyDB	MixingSecrets				Internal							
		L_a	L_{lr}	L_m	L_s	L_a	L_{lr}	L_m	L_s	L_a	L_{lr}	L_m	L_s
Base graph	50.7	1.45	1.42	198	7.30	2.16	2.02	22.9	1.12	.951	.940	1.63	
+ Gain/panning	g	.550	.583	.485	.550	.876	.856	.819	.973	.642	.619	.597	.734
+ Stereo imager	sg	.541	.564	.483	.553	.847	.834	.791	.928	.538	.616	.595	.727
+ Equalizer	e	sg	.450	.453	.390	.504	.700	.698	.622	.780	.522	.497	.467
+ Reverb	e	sg	r	.368	.361	.360	.390	.614	.601	.579	.674	.463	.451
+ Compressor	ec	sg	r	.315	.304	.297	.356	.558	.542	.512	.637	.396	.377
+ Noise gate	ecnsg	r	.302	.288	.281	.353	.548	.532	.502	.625	.393	.374	.343
+ Multitap delay (full)	ecnsgdr		.296	.288	.284	.324	.545	.529	.502	.618	.385	.369	.338
													.465

Table 2: Per-dataset results of the mixing consoles with different processor type configurations.

τ	MedleyDB							MixingSecrets							Internal												
	L_a	ρ	ρ_g	ρ_s	ρ_e	ρ_r	ρ_c	ρ_n	ρ_d	L_a	ρ	ρ_g	ρ_s	ρ_e	ρ_r	ρ_c	ρ_n	ρ_d	L_a	ρ	ρ_g	ρ_s	ρ_e	ρ_r	ρ_c	ρ_n	ρ_d
MC	—	.296	—	—	—	—	—	—	—	.545	—	—	—	—	—	—	—	—	.385	—	—	—	—	—	—	—	
BF	.01	.305	.63	.35	.81	.54	.77	.67	.71	.53	.566	.65	.50	.82	.43	.71	.61	.77	.75	.402	.80	.76	.92	.61	.81	.85	.87
DW	.01	.302	.56	.32	.79	.42	.74	.57	.56	.43	.561	.57	.47	.82	.22	.62	.54	.74	.55	.397	.74	.74	.82	.49	.70	.87	.86
H	.001	.295	.44	.22	.73	.26	.61	.50	.54	.24	.550	.43	.35	.76	.13	.42	.43	.55	.38	.388	.59	.48	.77	.40	.57	.78	.77
	.01	.302	.61	.32	.81	.48	.74	.69	.71	.52	.563	.62	.49	.85	.34	.64	.60	.79	.65	.400	.77	.73	.93	.56	.75	.85	.86
	.1	.375	.83	.59	.93	.83	.92	.80	.84	.90	.648	.84	.70	.91	.75	.86	.83	.93	.91	.474	.93	.90	.99	.84	.92	.93	.95

Table 3: Per-dataset results of the pruning. MC: mixing console. BF: brute-force, DW: dry/wet, and H: hybrid.

Algorithm 2 Music mixing graph search (dry/wet method).

Input: A mixing console G_c , dry tracks S , and mixture y

Output: Pruned graph G_p , parameters \mathbf{P} , and weights \mathbf{w}

```

1:  $\mathbf{P}, \mathbf{w} \leftarrow \text{Initialize}(G_c)$ 
2:  $\mathbf{P}, \mathbf{w} \leftarrow \text{Train}(G_c, \mathbf{P}, \mathbf{w}, S, y)$ 
3:  $L_a^{\min} \leftarrow \text{Evaluate}(G_c, \mathbf{P}, \mathbf{w}, S, y)$ 
4:  $G_p \leftarrow G_c$ 
5: for  $n \leftarrow 1$  to  $N_{\text{iter}}$  do
6:    $T_{\text{cand}} \leftarrow \text{GetProcessorTypeSet}(V)$ 
7:   for  $t$  in  $T_{\text{cand}}$  do
8:      $V_t, \mathbf{w}_t \leftarrow \text{Filter}(V, t), \text{Filter}(\mathbf{w}, t)$ 
9:      $N_t, r_t, \mathbf{m}_t \leftarrow |V_t|, 0.1, \mathbf{1}$ 
10:    end for
11:    while  $T_{\text{cand}} \neq \emptyset$  do
12:       $t \leftarrow \text{SampleType}(T_{\text{cand}})$ 
13:       $\bar{V}_t, \bar{\mathbf{m}} \leftarrow \text{GetLeastWeightNodes}(V_t, \mathbf{w}_t, \lfloor N_t r_t \rfloor)$ 
14:       $L_a \leftarrow \text{Evaluate}(G_p, \mathbf{P}, \mathbf{w} \odot \mathbf{m} \odot \bar{\mathbf{m}}, S, y)$ 
15:      if  $L_a < L_a^{\min} + \tau$  then
16:         $L_a^{\min} \leftarrow \min(L_a^{\min}, L_a)$ 
17:         $\mathbf{m} \leftarrow \mathbf{m} \odot \bar{\mathbf{m}}$ 
18:         $V_t \leftarrow V_t \setminus \bar{V}_t$ 
19:      else
20:        if  $N_t > 1$  then
21:           $r_t \leftarrow r_t / 2$ 
22:        else
23:           $T_{\text{cand}} \leftarrow T_{\text{cand}} \setminus \{t\}$ 
24:        end if
25:      end if
26:    end while
27:     $G_p, \mathbf{P}, \mathbf{w} \leftarrow \text{Prune}(G_p, \mathbf{P}, \mathbf{w}, \mathbf{m})$ 
28:     $\mathbf{P}, \mathbf{w} \leftarrow \text{Train}(G_p, \mathbf{P}, \mathbf{w}, S, y)$ 
29:  end for
30: return  $G_p, \mathbf{P}, \mathbf{w}$ 

```

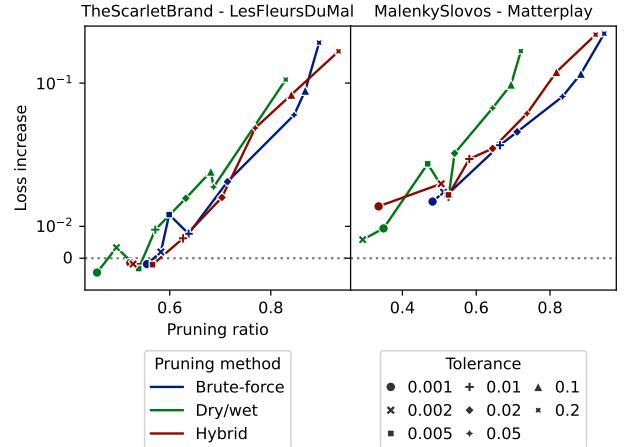


Figure 3: Loss increases from the mixing console and remaining processor ratios for different pruning methods and tolerances.

E. SUPPLEMENTARY RESULTS

- Table 2 and 3 report the per-dataset results on the mixing consoles and graph pruning, respectively.
- Figure 3 compares the pruning methods on 2 random-sampled songs using 7 tolerance settings from 0.001 to 0.2.
- Figure 4 shows multiple graphs obtained by pruning the same console (song) repeatedly.
- Refer to Figure 5-7 for more pruned graphs obtained with the default setting — hybrid method and $\tau = 0.01$.
- Figure 9-11 show more spectrogram plots.

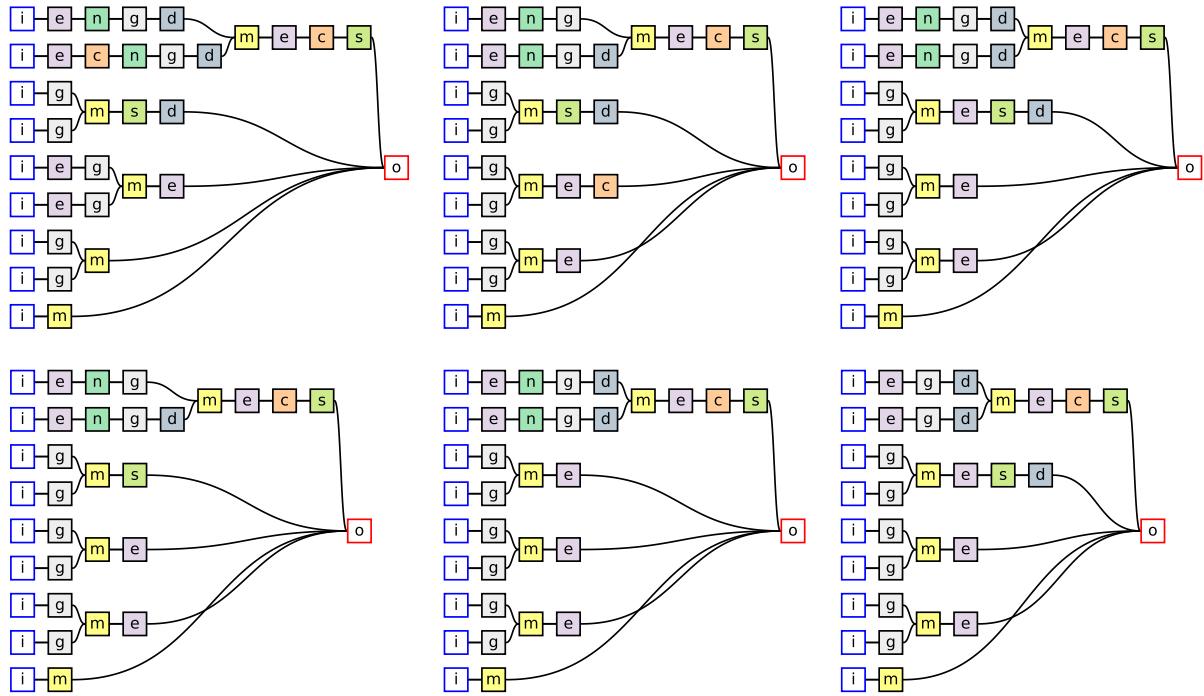


Figure 4: Each pruning run (default setting) yields a slightly different graph. Song: EthanHein_GirlOnABridge.

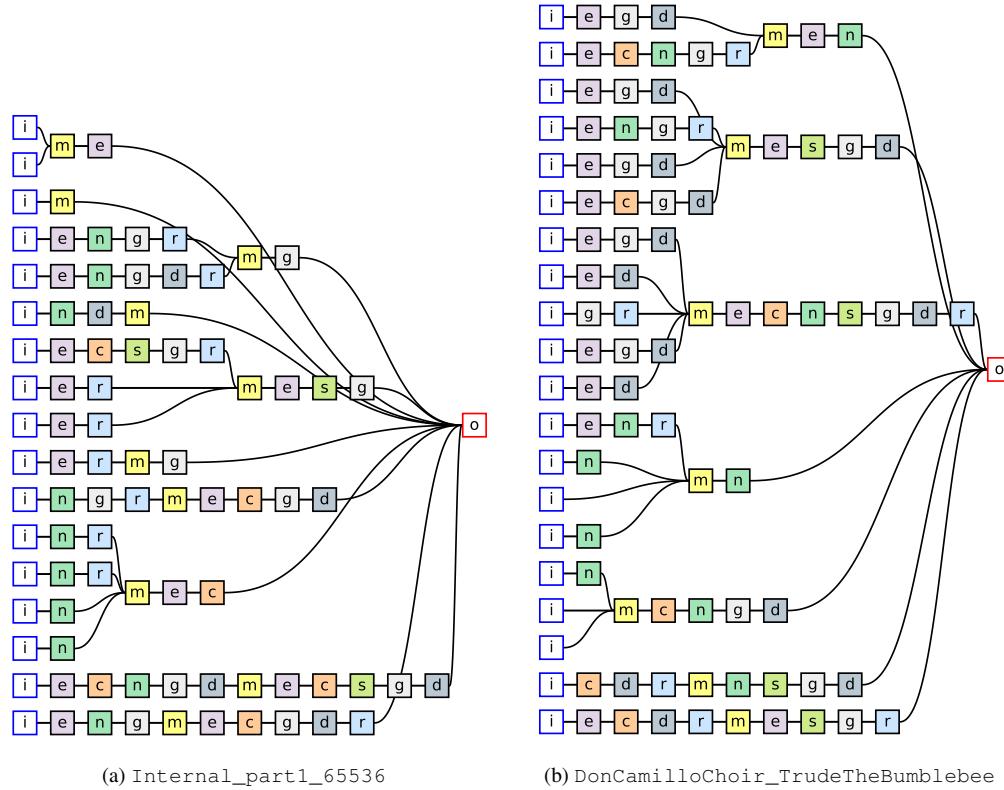


Figure 5: Example pruned graphs (default setting). the number of tracks: $K \leq 20$.

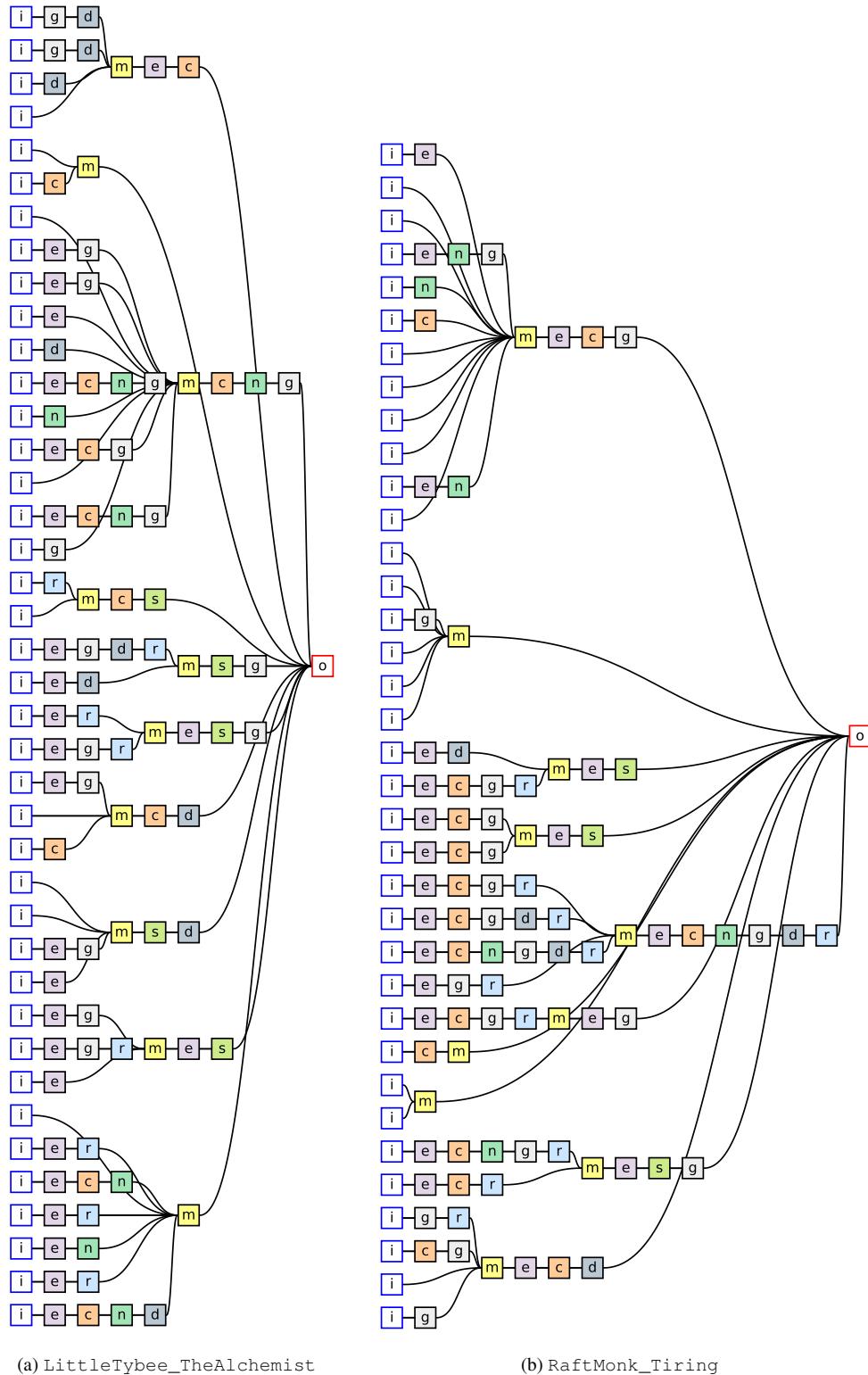


Figure 6: Example pruned graphs (default setting). the number of tracks: $K > 20$.

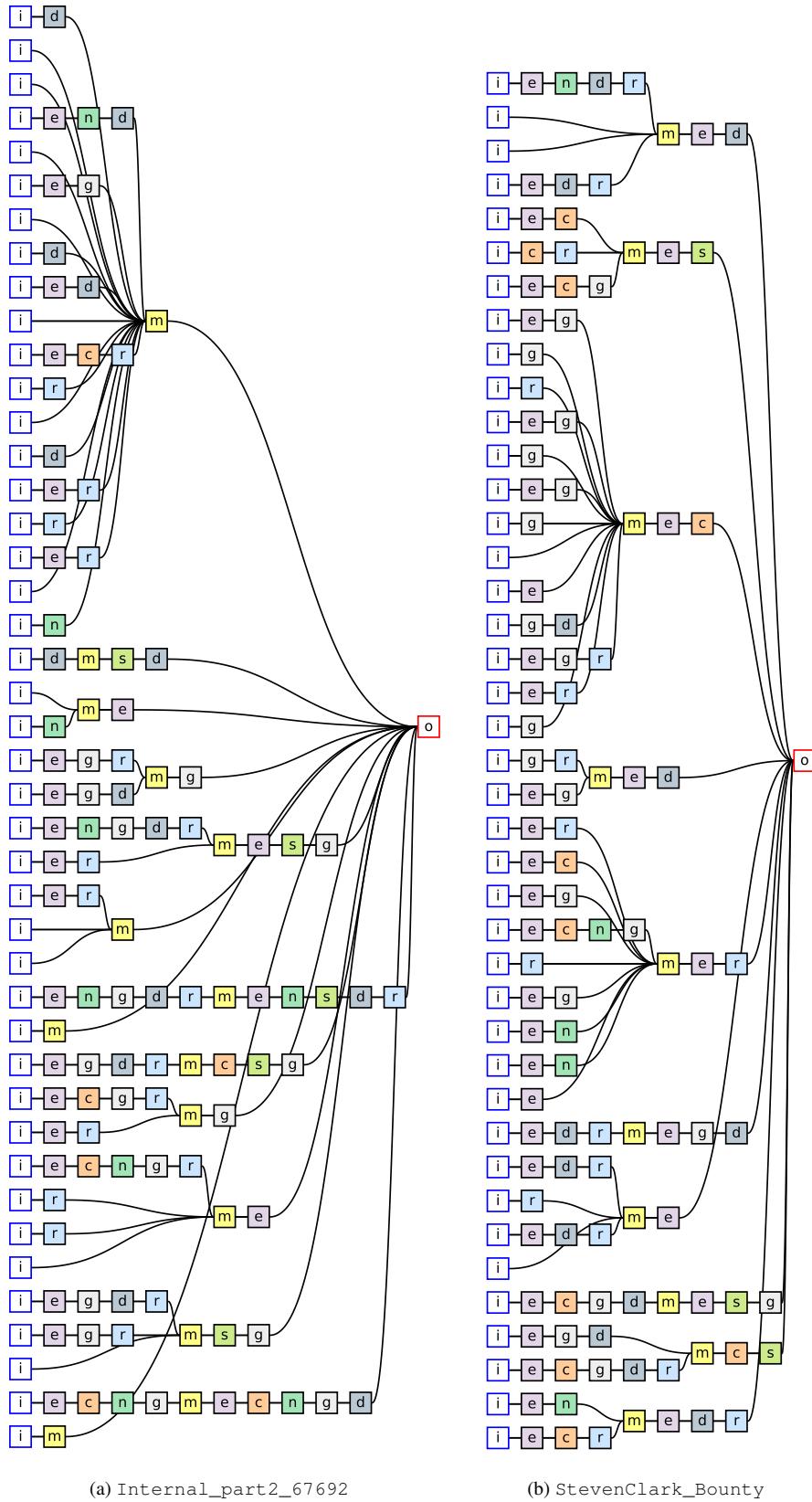


Figure 7: Example pruned graphs (default setting). the number of tracks: $K > 20$ (continued).

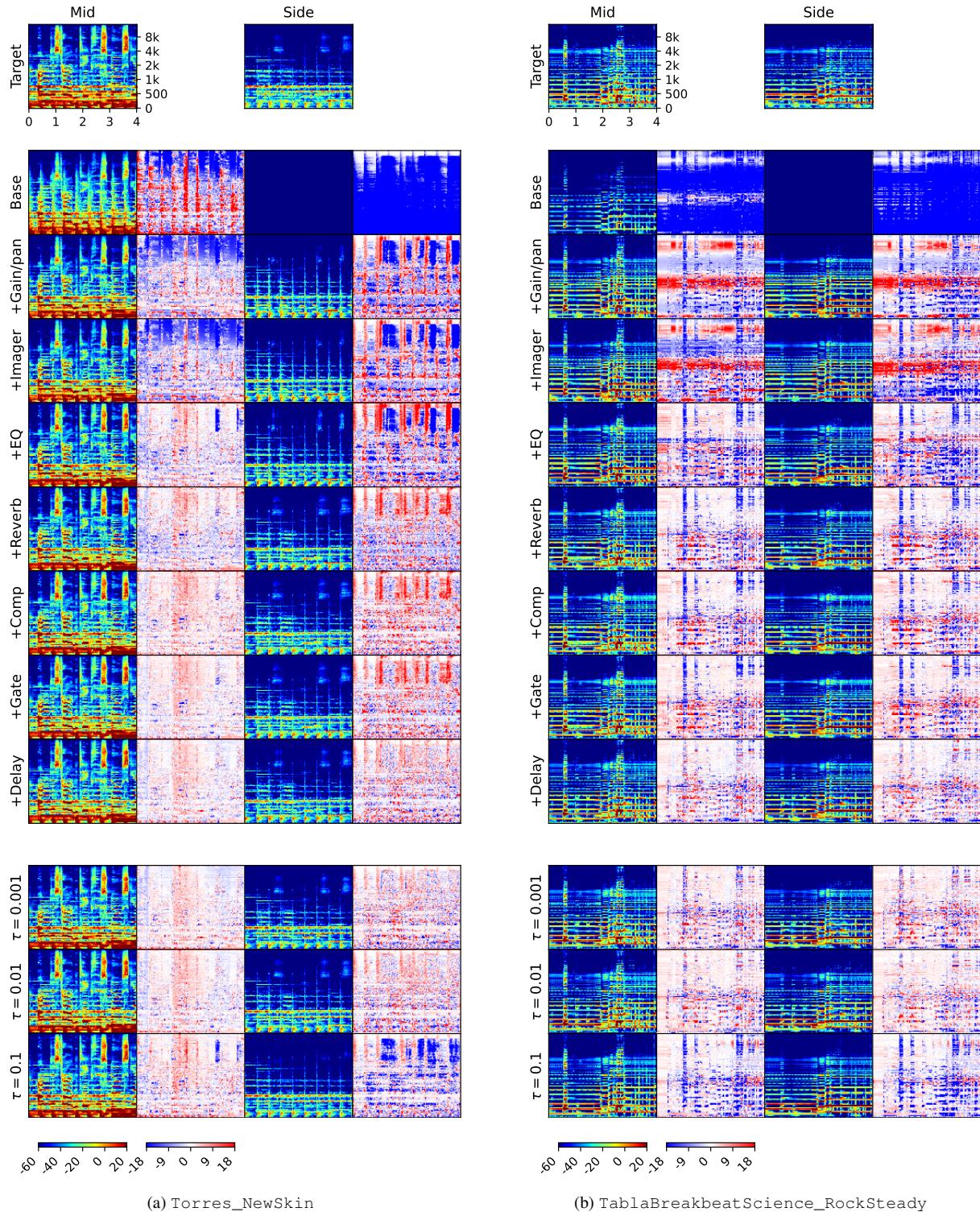


Figure 8: Matching of target music mixes with mixing consoles and their pruned versions: MedleyDB dataset.

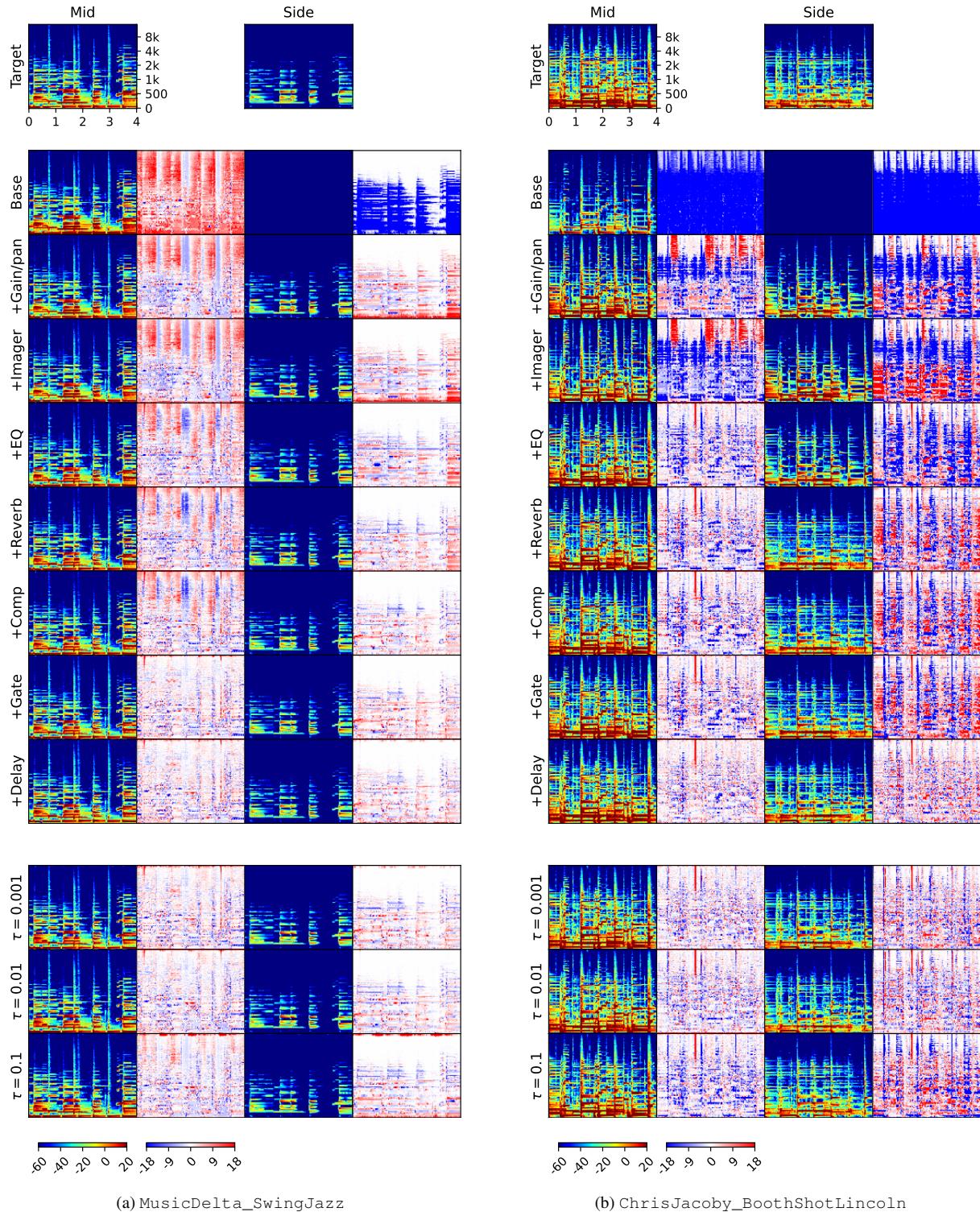


Figure 9: Matching of target music mixes with mixing consoles and their pruned versions: MedleyDB dataset.

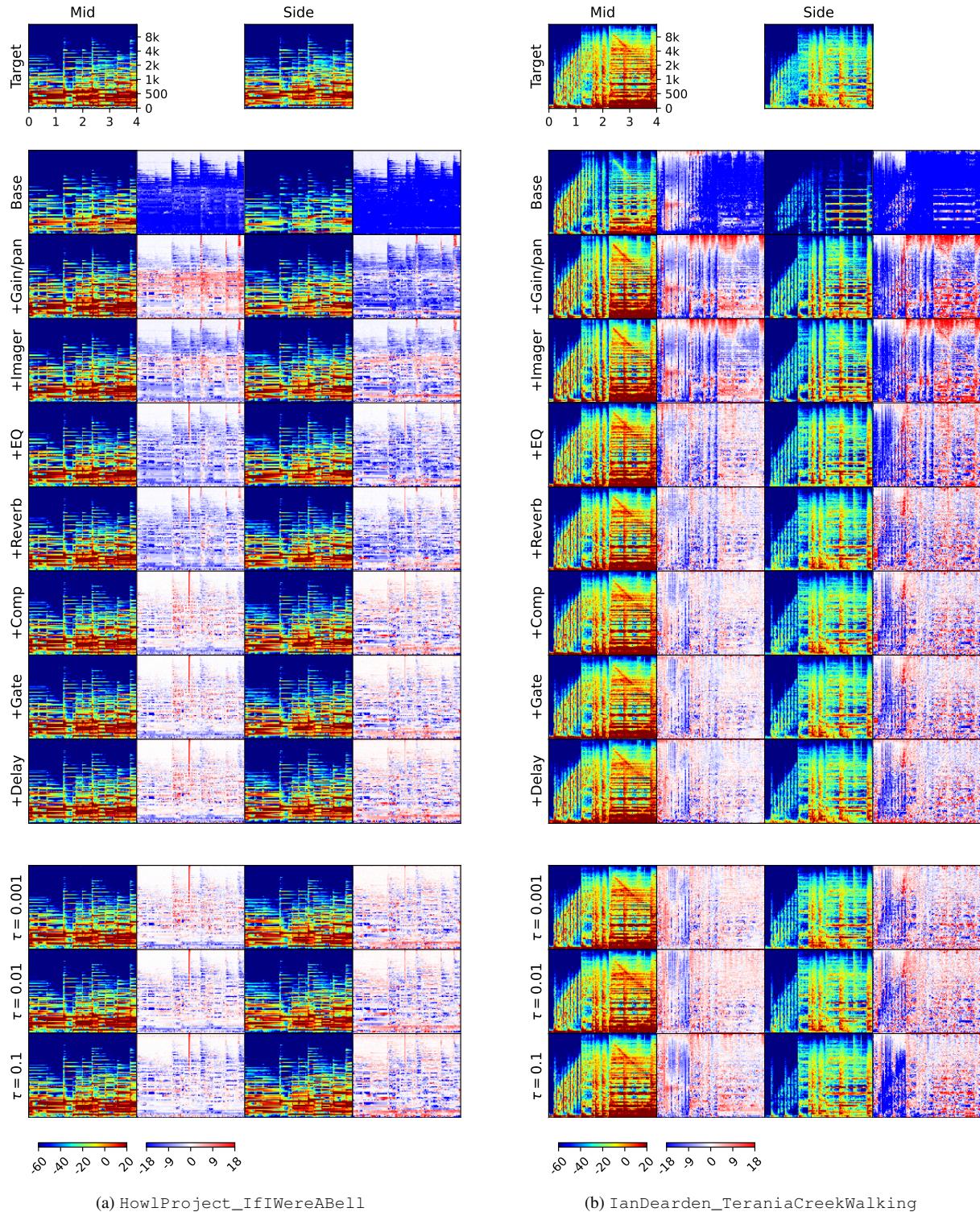


Figure 10: Matching of target music mixes with mixing consoles and their pruned versions: MixingSecrets dataset.

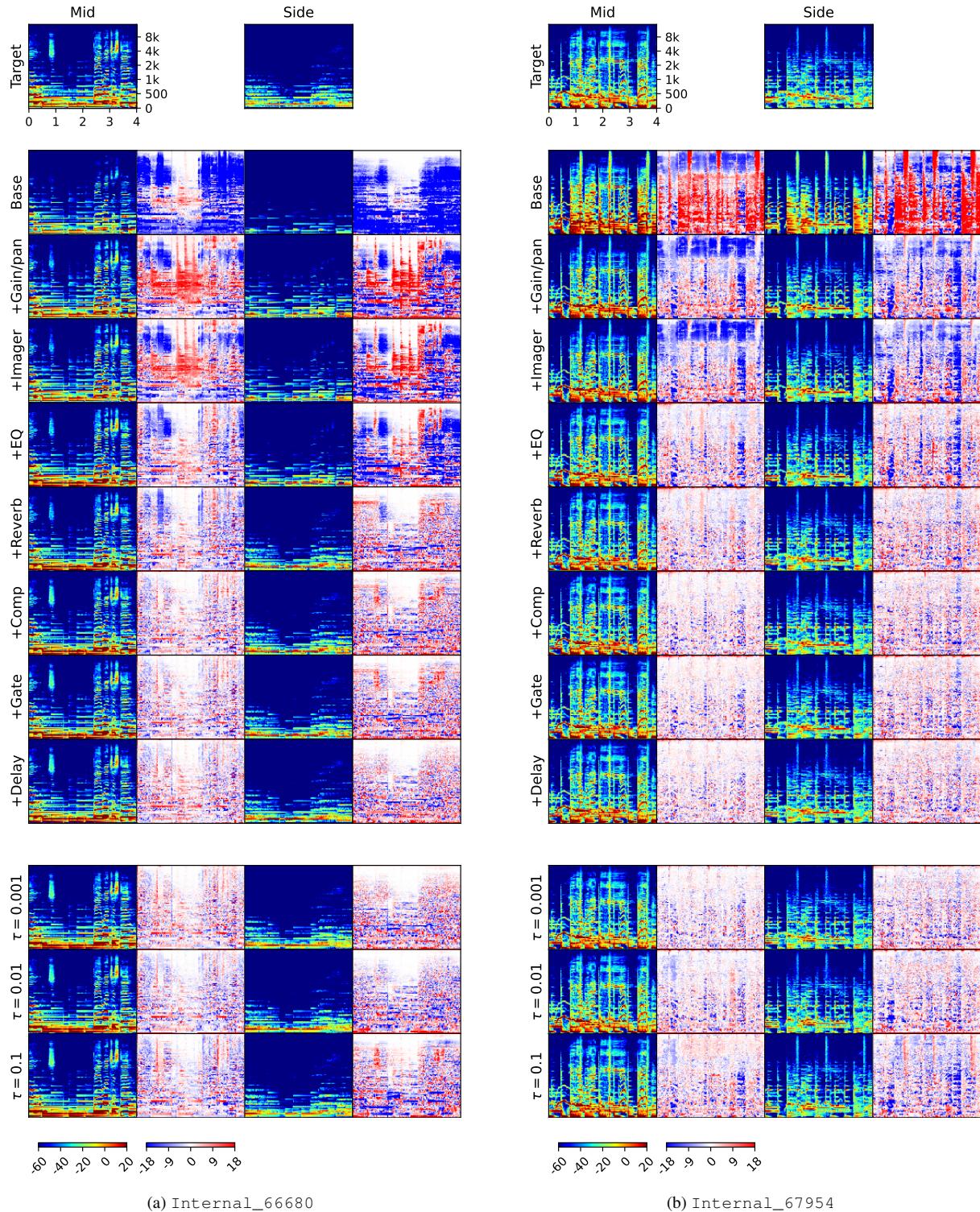


Figure 11: Matching of target music mixes with mixing consoles and their pruned versions: Internal dataset.

F. REFERENCES

- [1] S. Lee, J. Park, S. Paik, and K. Lee, “Blind estimation of audio processing graph,” in *IEEE ICASSP*, 2023.
- [2] C. Mitcheltree and H. Koike, “SerumRNN: Step by step audio VST effect programming,” in *Artificial Intelligence in Music, Sound, Art and Design*, 2021.
- [3] J. Guo and B. McFee, “Automatic recognition of cascaded guitar effects,” in *DAFx*, 2023.
- [4] J. Colonel, “Music production behaviour modelling,” 2023.
- [5] C. J. Steinmetz, J. Pons, S. Pascual, and J. Serrà, “Automatic multitrack mixing with a differentiable mixing console of neural audio effects,” in *IEEE ICASSP*, 2021.
- [6] N. Uzrad *et al.*, “DiffMoog: a differentiable modular synthesizer for sound matching,” *arXiv:2401.12570*, 2024.
- [7] J. Engel, L. H. Hantrakul, C. Gu, and A. Roberts, “DDSP: differentiable digital signal processing,” in *ICLR*, 2020.
- [8] M. A. Martínez-Ramírez, O. Wang, P. Smaragdis, and N. J. Bryan, “Differentiable signal processing with black-box audio effects,” in *IEEE ICASSP*, 2021.
- [9] Z. Ye, W. Xue, X. Tan, Q. Liu, and Y. Guo, “NAS-FM: Neural architecture search for tunable and interpretable sound synthesis based on frequency modulation,” *arXiv:2305.12868*, 2023.
- [10] F. Caspe, A. McPherson, and M. Sandler, “DDX7: Differentiable FM synthesis of musical instrument sounds,” in *ISMIR*, 2022.
- [11] A. Paszke *et al.*, “Pytorch: An imperative style, high-performance deep learning library,” *NeurIPS*, 2019.
- [12] C. J. Steinmetz, N. J. Bryan, and J. D. Reiss, “Style transfer of audio effects with differentiable signal processing,” *JAES*, vol. 70, no. 9, 2022.
- [13] J. T. Colonel, M. Comunità, and J. Reiss, “Reverse engineering memoryless distortion effects with differentiable wave-shapers,” in *AES Convention 153*, 2022.
- [14] A. Carson, S. King, C. V. Botinhao, and S. Bilbao, “Differentiable grey-box modelling of phaser effects using frame-based spectral processing,” in *DAFx*, 2023.
- [15] S. Nercessian, “Neural parametric equalizer matching using differentiable biquads,” in *DAFx*, 2020.
- [16] S. Lee, H.-S. Choi, and K. Lee, “Differentiable artificial reverberation,” *IEEE/ACM TASLP*, vol. 30, 2022.
- [17] B. Hayes, J. Shier, G. Fazekas, A. McPherson, and C. Saitis, “A review of differentiable digital signal processing for music & speech synthesis,” *Frontiers in Signal Process.*, 2023.
- [18] C. J. Steinmetz, T. Walther, and J. D. Reiss, “High-fidelity noise reduction with differentiable signal processing,” in *AES Convention 155*, 2023.
- [19] H. Liu, K. Simonyan, and Y. Yang, “DARTS: Differentiable architecture search,” in *ICLR*, 2019.
- [20] J. You, B. Liu, Z. Ying, V. Pande, and J. Leskovec, “Graph convolutional policy network for goal-directed molecular graph generation,” *NeurIPS*, 2018.
- [21] X. Serra and J. Smith, “Spectral modeling synthesis: A sound analysis/synthesis based on a deterministic plus stochastic decomposition,” *Computer Music Journal*, vol. 14, 1990.
- [22] J. T. Colonel and J. Reiss, “Reverse engineering a nonlinear mix of a multitrack recording,” *Journal of the AES*, vol. 71, no. 9, 2023.
- [23] D. Giannoulis, M. Massberg, and J. D. Reiss, “Digital dynamic range compressor design—a tutorial and analysis,” *JAES*, vol. 60, no. 6, 2012.
- [24] J. Colonel, J. D. Reiss, *et al.*, “Approximating ballistics in a differentiable dynamic range compressor,” in *AES Convention 153*, 2022.
- [25] B. Hayes, C. Saitis, and G. Fazekas, “Sinusoidal frequency estimation by gradient descent,” in *IEEE ICASSP*, 2023.
- [26] Y. Bengio, N. Léonard, and A. Courville, “Estimating or propagating gradients through stochastic neurons for conditional computation,” *arXiv:1308.3432*, 2013.
- [27] J. D. Ullman, “NP-complete scheduling problems,” *Journal of Computer and System Sciences*, vol. 10, no. 3, 1975.