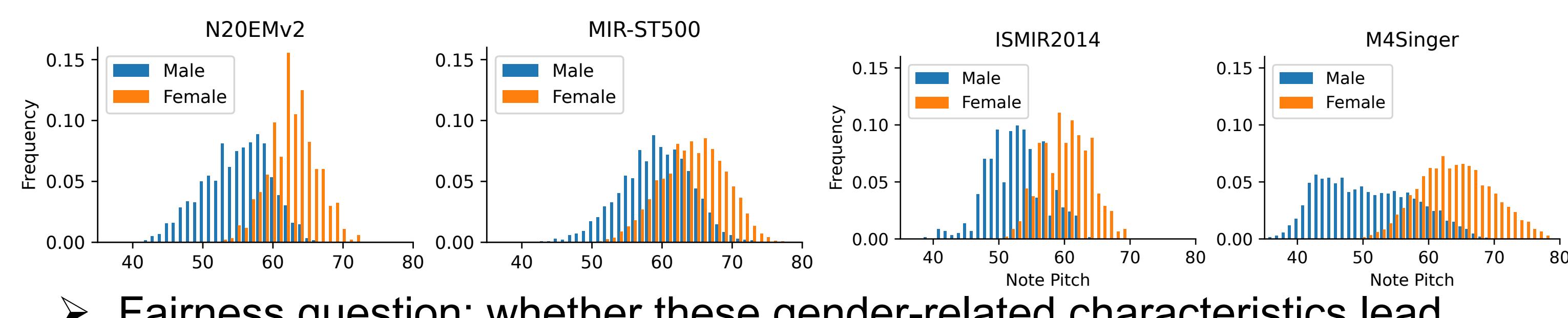


Elucidate Gender Fairness in Singing Voice Transcription

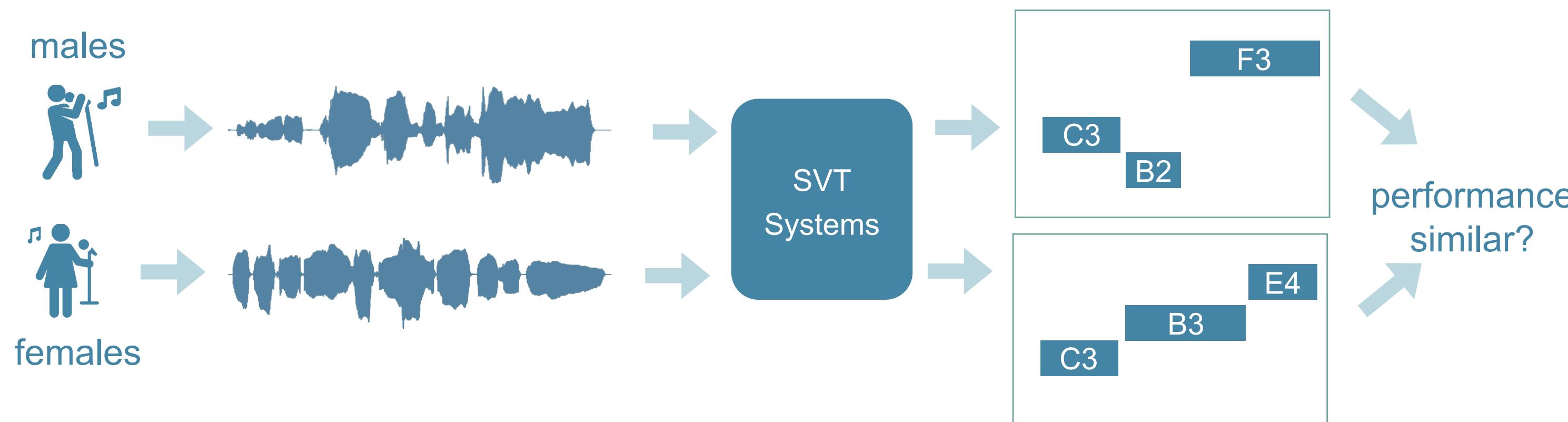
Xiangming Gu, Wei Zeng, Ye Wang

Motivation:

- Females and males sing differently. They have different sound characteristics, e.g. timbre and pitch distributions:



- Fairness question: whether these gender-related characteristics lead to a performance disparity in singing voice transcription.



- Why does fairness in SVT matter?

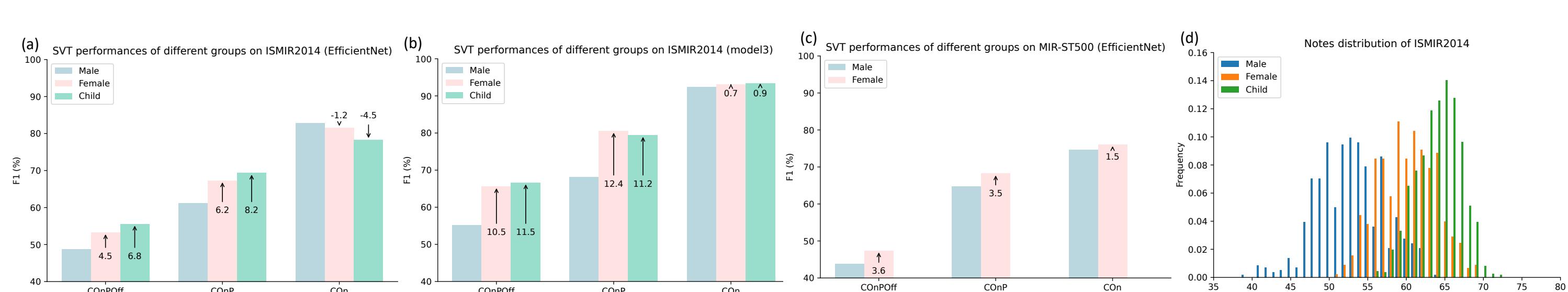
- Biased machine learning systems can lead to discriminatory treatment of certain groups.
- Biased SVT systems can lead to user inconvenience
- Biased SVT systems negatively affect downstream applications, e.g. singing voice synthesis

Fairness Analysis in SVT:

- Female superiority in SVT: evidence from four benchmark singing datasets

Dataset	Model	COnPOff (%)				CONP (%)				COn (%)			
		total ↑	female ↑	male ↑	gap	total ↑	female ↑	male ↑	gap	total ↑	female ↑	male ↑	gap
MIR-ST500	model1	52.39	53.11	51.47	-1.64	70.73	72.54	68.42	-4.12	78.32	79.36	77.00	-2.36
	model2	34.55	35.29	33.60	-2.55	51.64	52.76	50.21	-2.55	71.33	72.07	70.39	-1.68
	model3	52.84	54.02	51.33	-2.69	70.00	71.17	67.85	-3.32	78.05	78.84	77.05	-1.79
N20EMv2	model1	55.20	60.96	51.55	-9.41	72.03	79.76	67.11	-12.65	88.51	90.63	87.16	-3.47
	model2	68.62	74.82	64.68	-10.14	75.69	81.27	72.14	-9.12	92.83	94.33	91.88	-2.44
	model3	73.06	78.34	69.69	-8.65	79.56	84.38	76.49	-7.89	93.66	94.96	92.83	-2.13
ISMIR2014	model1	52.58	61.22	45.28	-15.94	67.75	74.90	61.70	-13.20	92.13	91.93	92.30	+0.37
	model2	57.35	62.42	53.06	-9.36	72.15	79.35	66.06	-13.29	91.53	92.25	90.92	-1.33
	model3	59.95	65.61	55.16	-10.45	73.85	80.55	68.19	-12.36	92.80	93.18	92.49	-0.69
M4Singer	wav2vec2	53.66	57.27	49.93	-7.34	61.95	66.36	57.38	-8.98	82.60	81.66	83.58	+1.91
	Hubert	55.11	58.39	51.73	-6.66	64.17	68.10	60.10	-8.00	82.13	81.53	82.76	+1.23
	wavLM	57.06	60.33	53.68	-6.66	65.40	68.87	61.82	-7.05	82.73	82.29	83.19	+0.90
	data2vec	53.98	57.45	50.40	-7.05	63.04	67.09	58.86	-8.23	82.30	82.13	82.48	+0.35

- Female superiority in SVT: evidence from SSL-based SVT models and EfficientNet-based SVT models



- Potential reasons:

- Rule out the gender data imbalance
- Hypothesis: pitch distributions make substantial contributions.

Fairness formulation in SVT:

- Training data

$$\mathcal{D}_S = \{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)}, A^{(n)})\}_{n=1}^{N_S} \quad \mathcal{D}_{\mathcal{T}} = \{(\mathbf{x}^{(n)}, \mathbf{y}^{(n)}, A^{(n)})\}_{n=1}^{N_{\mathcal{T}}}$$

- SVT metric

COnPOff: correct onset, offset, pitch

- Utility metric

$$U = \mathbb{E}_{(\mathbf{x}, \mathbf{y}) \sim \mathcal{D}_{\mathcal{T}}} [\text{metric}(\hat{\mathbf{y}}, \mathbf{y})]$$

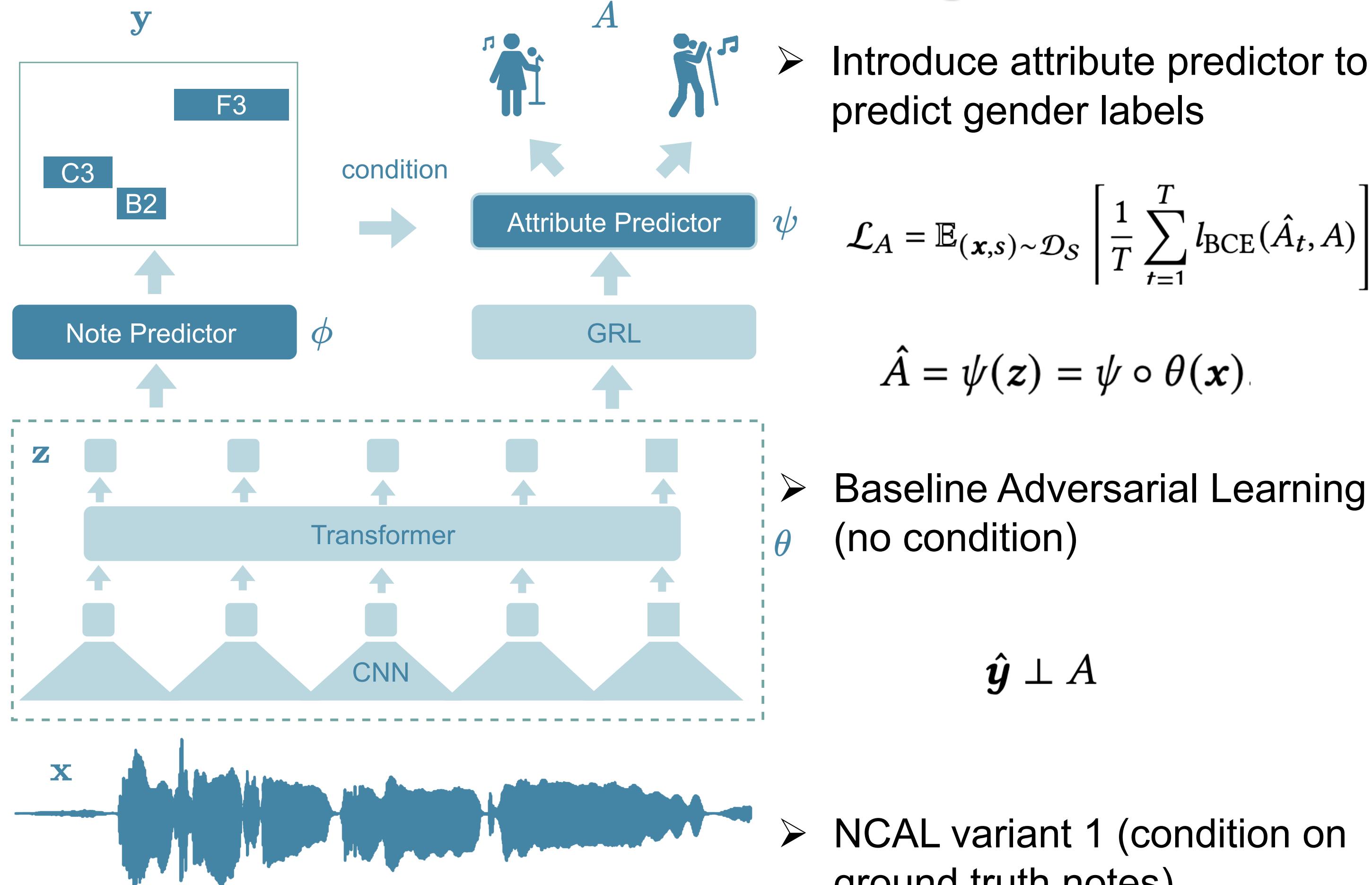
- Fairness metric ($A = 1$: males $A = 0$: females)

$$F = \mathbb{E}_{(\mathbf{x}, \mathbf{y}, A=1) \sim \mathcal{D}_{\mathcal{T}}} [\text{metric}(\hat{\mathbf{y}}, \mathbf{y})] - \mathbb{E}_{(\mathbf{x}, \mathbf{y}, A=0) \sim \mathcal{D}_{\mathcal{T}}} [\text{metric}(\hat{\mathbf{y}}, \mathbf{y})]$$

- Fairness-Utility tradeoff

$$\max \min\{F, 0\}, \quad \max U$$

Note-conditioned Adversarial Learning for Fairness:



Introduce attribute predictor to predict gender labels

$$\mathcal{L}_A = \mathbb{E}_{(\mathbf{x}, s) \sim \mathcal{D}_S} \left[\frac{1}{T} \sum_{t=1}^T l_{BCE}(\hat{A}_t, A) \right]$$

$$\hat{A} = \psi(z) = \psi \circ \theta(\mathbf{x}).$$

Baseline Adversarial Learning (no condition)

$$\hat{y} \perp A$$

NCAL variant 1 (condition on ground truth notes)

$$z \perp A | y$$

NCAL variant 2 (condition on note predictions): align different pitches simultaneously

$$z \perp A | \hat{y}$$

Algorithm 1 Note-conditioned adversarial learning for SVT

Require: Acoustic encoder $\theta^{(0)}$ pre-trained on speech data under SSL objective, randomly initialized label predictor $\phi^{(0)}$ and attribute predictor $\psi^{(0)}$, learning rates η_1, η_2, η_3 for θ, ϕ, ψ , training steps K_1, K_2 for linear probing and full finetuning.

```

for k = 1 to  $K_1 + K_2$  do
   $z = \theta^{(k-1)}(\mathbf{x}), \hat{y} = \phi^{(k-1)}(z)$ 
   $\hat{A}_t = \psi^{(k-1)}(z.\text{detach}(), \hat{y}.\text{detach}()),$  compute  $\mathcal{L}_A$  in Eq. 6
   $\psi^{(k)} = \psi^{(k-1)} - \eta_3 \frac{\partial \mathcal{L}_A}{\partial \psi^{(k-1)}}$   $\triangleright$  Update  $\psi$ 
   $\hat{A}_t = \psi^{(k)}(z, \hat{y}),$  compute  $\mathcal{L}_y, \mathcal{L}_A$  in Eq. 2 and 6
   $\phi^{(k)} = \phi^{(k-1)} - \eta_2 \frac{\partial \mathcal{L}_y}{\partial \phi^{(k-1)}}$   $\triangleright$  Update  $\phi$ 
  if  $k \leq K_1$  then
     $\theta^{(k)} = \theta^{(k-1)}$   $\triangleright$  Update  $\theta$ 
  else
     $\theta^{(k)} = \theta^{(k-1)} - \eta_1 (\frac{\partial \mathcal{L}_y}{\partial \theta^{(k-1)}} - \lambda \frac{\partial \mathcal{L}_A}{\partial \theta^{(k-1)}})$ 
  end if
end for

```

Enforce Fairness in SVT:

- Gender Bias Mitigation Results of NCAL.

Train set	Test set	Method	COnPOff (%)		CONP (%)	
			Utility (U) ↑	Fairness (F) ↑	Utility (U) ↑	Fairness (F) ↑
M4Singer	M4singer	ERM	53.66	- 7.34	61.95	- 8.98
		Ours	52.48 (-1.18)	- 3.61 (+3.73)	60.67 (-1.28)	- 4.21 (+4.77)
MIR-ST500	N20EMv2	ERM	55.20	- 9.41	72.03	- 12.65
		Ours	53.29 (-1.91)	- 4.14 (+5.27)	72.71 (+0.68)	- 10.82 (+1.83)
MIR-ST500	ISMIR2014	ERM	52.58	- 15.94	67.75	- 13.20
		Ours	48.18 (-4.40)	- 8.29 (+7.65)	65.40 (-2.35)	- 9.08 (+4.12)
MIR-ST500	N20EMv2	ERM	73.06	- 8.65	79.56	- 7.89
		Ours	72.43 (-0.63)	- 5.78 (+2.87)	78.47 (-1.09)	- 6.82 (+1.07)
N20EMv2	ISMIR2014	ERM	59.95	- 10.45	73.85	- 12.36
		Ours	59.57 (-0.38)	- 7.49 (+2.96)	73.33 (-0.52)	- 7.75 (+4.61)

- Why do we need fairness-utility tradeoff?

Train set	Test set	Method	η_3	λ	COnPOff (%)		CONP (%)	
Utility (U) ↑	Fairness (F) ↑	Utility (U) ↑	Fairness (F) ↑					

<tbl_r cells="8" ix="5" maxcspan="1" maxrspan="1" usedcols