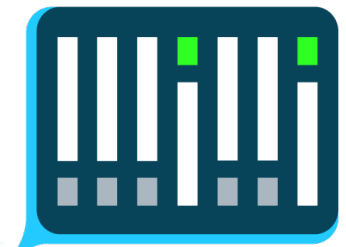




The
University
Of
Sheffield.



Exploring the Use of Group Delay for Generalised VTS-based Noise Compensation

Erfan Loweimi,

Jon Barker and Thomas Hain

{e.loweimi, j.p.barker, t.hain}@sheffield.ac.uk



Presented by:
Yannis Stylianou

Speech and Hearing Research Group (**SPANDH**)



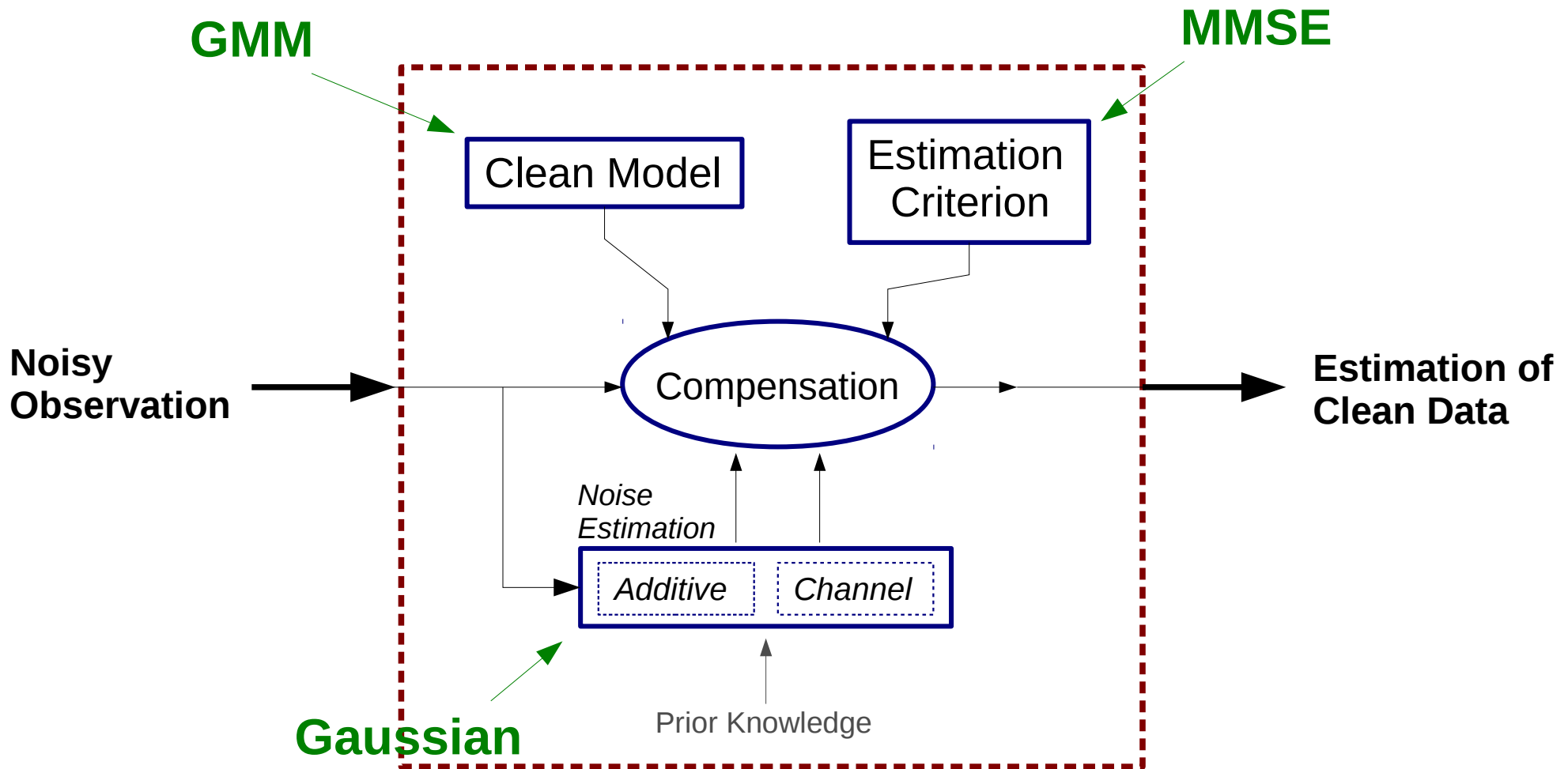
Outline

- **Generalised VTS** (gVTS) approach to Robust ASR
- Extension of the gVTS to group delay domain
 - Environment Model
 - Challenges
 - Proposed Solutions
 - Deriving Equations
- Experimental Results and Discussion





Model-based Noise Compensation



(g)VTS is a model-based technique for noise compensation.





(g)VTS Pseudocode

0. GMM of **CLEAN**

– *For each utterance ...*

1. Compute the environment model
2. Apply the (Gen)Log
3. Factor out **CLEAN** part and compute the distortion function
4. Estimate **Noise**
 - 4.1. Additive
 - 4.2. Channel
5. Linearise using **Taylor series**
 - 5.1. Points → means of Gaussians
 - 5.2. Jacobians → partial derivatives
6. Estimate **CLEAN** features using **MMSE**





Advantages of gVTS

- gVTS → replacing *Log* with *GenLog* in VTS
- One extra degree of freedom (α)
 - A non-linear transform with **statistical** effect [App. 1]
 - Can improve linearity, homoschodasticity and Gaussianty
 - Compensation is carried out in a space with a higher signal-to-noise ratio (SNR) [App. 2]
 - Further **robustness**
 - The optimal value for α → 0.05 – 0.1 [App. 3]

$$\begin{cases} GenLog(x; \alpha) = \frac{1}{\alpha}(x^\alpha - 1), & x > 0 \quad \alpha \neq 0 \\ \lim_{\alpha \rightarrow 0} GenLog(x; \alpha) = \log(x) \end{cases}$$



Extension of the gVTS To Group Delay (GD) Domain

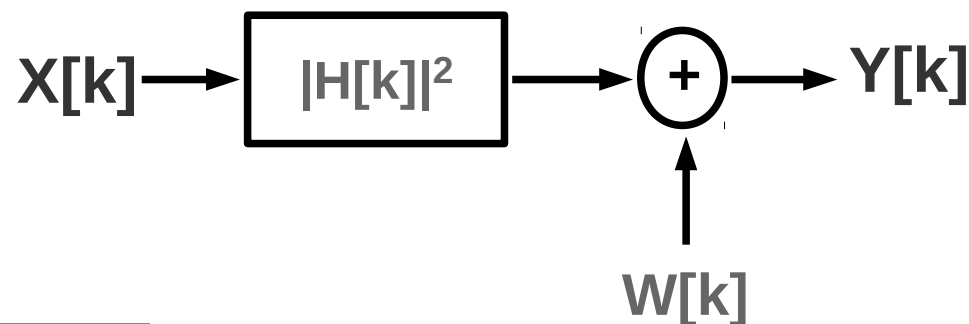




Environment Model

Periodogram domain

$$|Y|^2 = |X|^2 |H|^2 + |W|^2$$



Group delay domain

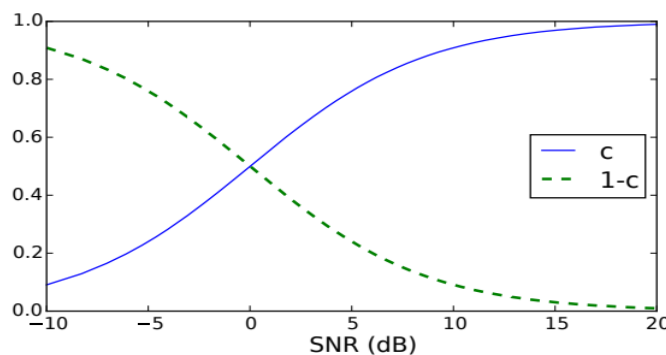
$$\tau_Y = \frac{|X|^2 |H|^2}{|Y|^2} (\tau_X + \tau_H) + \frac{|W|^2}{|Y|^2} \tau_W$$

$$\tau_Y = \frac{\xi}{1 + \xi} \tau_X + \frac{1}{1 + \xi} \tau_W = c \tau_X + (1 - c) \tau_W$$

convex combination
←

$$\xi = \frac{|X|^2}{|W|^2} : \text{a priori SNR}$$

$$c = \frac{\xi}{1 + \xi}$$





Challenges

1) Larger number of variables

- **4** in periodogram domain vs **8** in GD domain
 - For each variable a statistical model should be estimated
 - Noise compensation would be more complicated

2) Dynamic range compression using log and/or power transformation is problematic

- Group delay can be negative

4 variables \longrightarrow $|Y|^2 = |X|^2 |H|^2 + |W|^2$

8 variables \longrightarrow $\tau_Y = \frac{|X|^2 |H|^2}{|Y|^2} (\tau_X + \tau_H) + \frac{|W|^2}{|Y|^2} \tau_W$



Larger Number of Variables

- How to reduce number of variables?
 - Variables representing similar information and are added/multiplied may be encapsulated into one variable, e.g. group delay and power spectrum
 - Variables tends to zero in expected sense, may be removed, e.g. clean signal and noise cross-correlation

$$|Y|^2 \tau_Y = |X|^2 |H|^2 (\tau_X + \tau_H) + |W|^2 \tau_W$$

$$Q_Y = Q_X |H|^2 + Q_H |X|^2 + Q_W$$

Group delay-power
Product spectrum (PS)

#variables: 6
Still larger than periodogram
domain which is 4!



Larger Number of Variables ...

$$|Y|^2 = |X|^2 |H|^2 + |W|^2$$

← Periodogram domain

$$Q_Y = Q_X |H|^2 + \overbrace{\tau_H |H|^2 |X|^2}^{Q_H} + Q_W$$

← Product Spectrum domain

Extra *undesired* term!

- Without it the environment model resembles the periodogram domain
- If $\tau_H |H|^2 |X|^2$ tends to zero, (g)VTS equations in periodogram domain can be used
- Obviously $|X|^2$ and $|H|^2$ are not zero!
- What about τ_H ?



Phase and Group Delay of Channel

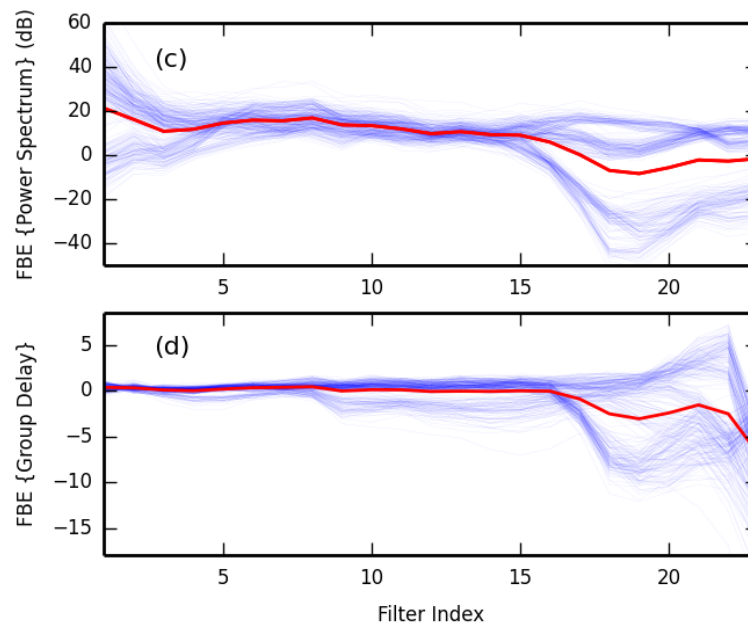
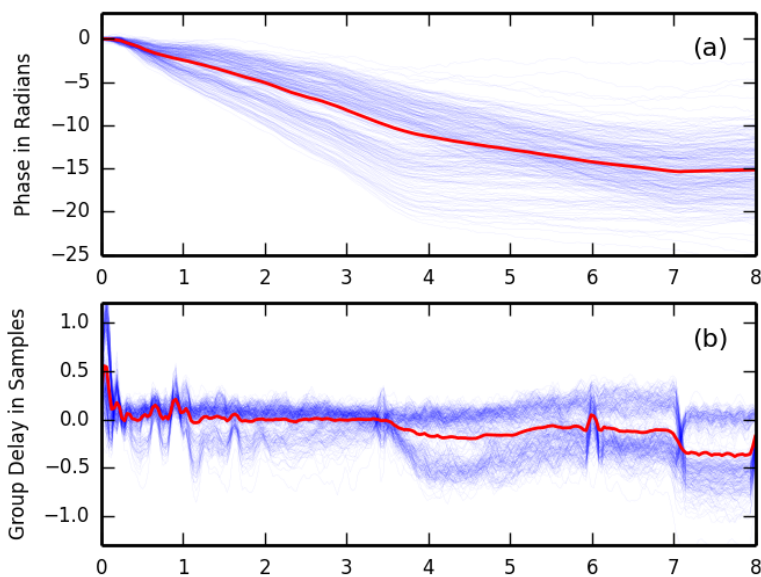
- Test sets A and C of Aurora-4 database may be used as *stereo* data to estimate channel Fourier transform

$$\begin{cases} \text{Test Set A} \Rightarrow Y^A = X \\ \text{Test Set C} \Rightarrow Y^C = X H \end{cases} \Rightarrow H_t = \frac{Y_t^C}{Y_t^A}$$

t: Frame index



Channel phase spectrum



Filterbank Energies (FBE) – Periodogram

(FBE) – group delay

Red curve: average over 330 utterances



Channel Phase Spectrum and Group Delay

- In the expected sense, the group delay of the channel tends to zero, as a result ...
 - *Undesired* term can be removed
 - Equation would be similar to periodogram domain

$$|Y|^2 = |X|^2 |H|^2 + |W|^2$$

$$Q_Y = Q_X |H|^2 + \tau_H |H|^2 |X|^2 + Q_W$$

Note: A red dashed box highlights the term $\tau_H |H|^2 |X|^2$, with a red arrow pointing to a red '0' above it, indicating it is to be removed.

 $Q_Y \approx Q_X |H|^2 + Q_W$



Dynamic Ranges Compression

- Dynamic range of product spectrum is comparable to power spectrum
 - Needs compression before statistical modelling
- (Gen)Log cannot be applied directly because group delay can be negative
- Possible solutions:

Absolute value \rightarrow $|Q_Z|^\alpha$

$sign(Q_Z)|Q_Z|^\alpha$ \leftarrow Use sign function and abs

Add a constant \rightarrow $(Q_Z + C)^\alpha$

Best solution

$\check{Q}_Z = max(Q_Z, 0)^\alpha$ \leftarrow Flooring



gVTS in the GD-power Product Spectrum Domain

1. Statistical models

$$\begin{cases} \check{Q}_X \sim \sum_{m=1}^M p_m^{\check{Q}_X} \mathcal{N}(\mu_m^{\check{Q}_X}, \Sigma_m^{\check{Q}_X}) \\ \check{Q}_W \sim \mathcal{N}(\mu^{\check{Q}_W}, \Sigma^{\check{Q}_W}) \\ \check{H} \sim \mathcal{N}(\mu^{\check{H}}, \Sigma^{\check{H}}) \end{cases}$$

2. Environment model (after applying GenLog)

$$\check{Q}_Y \approx \check{Q}_X \check{H} \underbrace{\left(1 + \left(\frac{\check{Q}_W}{\check{Q}_X \check{H}} \right)^{\frac{1}{\alpha}} \right)^\alpha}_{\text{Distortion function } \rightarrow \check{G}(\check{Q}_X, \check{Q}_W, \check{H})}$$

3. Taylor series (linearisation)

$$\begin{aligned} \check{Q}_Y \approx & \check{Q}_{Y0} + J^{\check{Q}_X} (\check{Q}_X - \check{Q}_{X0}) + \\ & + J^{\check{Q}_W} (\check{Q}_W - \check{Q}_{W0}) + J^{\check{H}} (\check{H} - \check{H}_0) \end{aligned}$$





gVTS in the GD-power Product Spectrum Domain

4. Compute the Jacobians (partial derivatives)

$$\check{V}_m = \left(\frac{\mu^{\check{Q}_w}}{\mu_m^{\check{Q}_x} \mu^{\check{H}}} \right)^{\frac{1}{\alpha}} \longleftrightarrow$$

$$J_m^{\check{Q}_x} = \frac{\partial \check{Q}_Y}{\partial \check{Q}_x} = \text{diag} \left\{ \frac{\mu^{\check{H}}}{(1 + \check{V}_m)^{1-\alpha}} \right\}$$

$$J_m^{\check{Q}_w} = \frac{\partial \check{Q}_Y}{\partial \check{Q}_w} = \text{diag} \left\{ \left(\frac{\check{V}_m}{1 + \check{V}_m} \right)^{1-\alpha} \right\}$$

$$J_m^{\check{H}} = \frac{\partial \check{Q}_Y}{\partial \check{H}} = \text{diag} \left\{ \frac{\mu_m^{\check{Q}_x}}{(1 + \check{V}_m)^{1-\alpha}} \right\}$$

5. Compute noisy observations (Q_Y) statistics

$$\check{Q}_Y \sim \sum_{m=1}^M p_m^{\check{Q}_Y} \mathcal{N}(\mu_m^{\check{Q}_Y}, \Sigma_m^{\check{Q}_Y}) \longleftrightarrow$$

$$p_m^{\check{Q}_Y} \approx p_m^{\check{Q}_x}$$

$$\mu_m^{\check{Q}_Y} \approx \mu_m^{\check{Q}_x} \mu^{\check{H}} \left(1 + \left(\frac{\mu^{\check{Q}_w}}{\mu_m^{\check{Q}_x} \mu^{\check{H}}} \right)^{\frac{1}{\alpha}} \right)^\alpha$$

$$\Sigma_m^{\check{Q}_Y} \approx J_m^{\check{Q}_x} \Sigma_m^{\check{Q}_x} J_m^{\check{Q}_x T} + J_m^{\check{Q}_w} \Sigma^{\check{Q}_w} J_m^{\check{Q}_w T} + J_m^{\check{H}} \Sigma^{\check{H}} J_m^{\check{H} T}$$





gVTS in the GD-power Product Spectrum Domain

6. MMSE estimate

$$\hat{Q}_X^{MMSE} = \mathbb{E}[\check{Q}_X | \check{Q}_Y] = \int \check{Q}_X p(\check{Q}_X | \check{Q}_Y) d\check{Q}_X$$
$$\approx \check{Q}_Y \sum_{m=1}^M p(m | \check{Q}_Y) \frac{1}{\check{G}(\mu_m^{\check{Q}_X}, \mu^{\check{Q}_W}, \mu^{\check{H}})}$$



Experimental Setup

- Database: Aurora-4
- Training sets: (each set: 7138 utterances, ~ 14 hours)
 - Clean-SI-84 → only clean data → **CL**
 - Noisy-SI-84 → clean and additive noise, SNR: 15 dB → **M1**
 - Multi-SI-84 → clean+additive+channel, SNR: 15 dB → **M2**
- Test set: Eval-92 → 330 utterances, ~ 40 minutes
 - 14 noise types artificially added using FaNT tool → 4620 utterances, grouped into
 - Test set A: Clean
 - Test set B: Additive noise, SNR: 10 dB
 - 6 noise types: Airport, Babble, Car, Restaurant, Street, Train Station
 - Test set C: Channel distortion
 - Test set D: Additive and Channel noise, SNR: 10 dB (6 additive noise types+channel distortion)
- Channel estimation → Method proposed in our earlier publication [App. 4]
- GMM/HMM → HTK → state-clustered triphones → 16 Gaussians, 4 iterations
- DNN (TNET) → 4 hidden layers (1300 nodes each) → bottleneck (26 nodes) → output-layer [App. 5]
- Language model → bigram (perplexity: 147)



Experimental Results (WER)

Aurora4 – GMM/HMM

Feature	α	A	B	C	D	Ave
MFCC-CL	log	7.0	33.7	23.6	49.9	28.6
MFCC-M1	log	9.1	18.4	23.4	35.9	21.7
MFCC-M2	log	10.7	17.0	19.1	31.3	19.5
PS-CL	log	7.1	33.7	23.7	49.9	28.6
gPS-CL	0.05	7.0	25.3	23.2	42.9	24.6
gPS-CL	0.1	8.1	22.1	25.6	40.8	24.1
gVTS-CL	0.05	6.5	20.2	13.9	34.3	18.7
gVTS-CL	0.075	7.1	19.8	15.0	34.0	19.0
gVTS-CL	0.1	7.4	19.6	15.4	33.9	19.1

Relative WER reduction (relative to CL):

A → 7.7% C → 41.1% Ave → 34.8%

B → 41.8% D → 32.1%

Relative to M1:

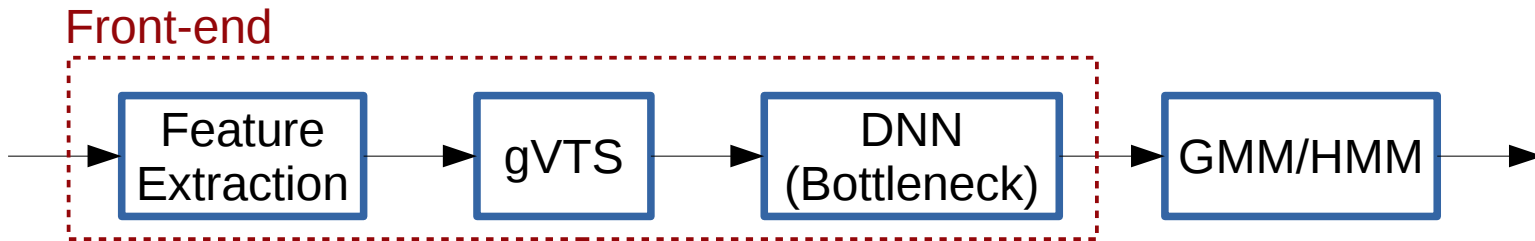
Ave → +13.8%

Relative to M2:

Ave → +4.1%

Experimental Results (WER)

Aurora4 – Bottleneck



Feature	α	A	B	C	D	<i>Ave</i>
BN{gPS}-CL	0.1	5.5	24.2	26.8	45.4	25.5
BN{gVTS}-CL	0.1	4.6	20.6	16.0	36.7	19.5
BN{gPS}-M1	0.1	5.5	11.1	23.5	32.3	18.1
BN{gVTS}-M1	0.1	5.3	12.4	14.3	30.6	15.6
BN{gPS}-M2	0.1	5.7	10.8	13.0	24.7	13.6
BN{gVTS}-M2	0.1	5.6	11.9	12.3	26.5	14.1

Relative WER reduction (relative to CL):
 A → 16.4% C → 40.3% **Ave** → 23.5%
 B → 18.2% D → 24.9%

Relative to M1:
Ave → +13.8%

Relative to M2:
Ave → -3.7%



Discussion and Conclusion

- gVTS and channel estimation techniques, proposed in earlier publications, successfully extended to product spectrum domain
- On average, the propose system trained by only clean data outperforms the MFCC-based system trained by Multi-style data (**M2**) in conventional GMM/HMM
- Combination of the gVTS and DNN is
 - **Super-additive** [App. 6] when there is a structural mismatch between the test and train conditions, e.g. **CL** or **M1** training conditions
 - Allows for building a robust system using DNN even when only clean data is available
 - **sub-additive** [App. 6] when all noise types with comparable SNR are available during training (**M2**)
 - In this case, **discrimination** is the main issue, not **robustness**





That's it!

- Thanks for your attention
- Q&A





Appendices

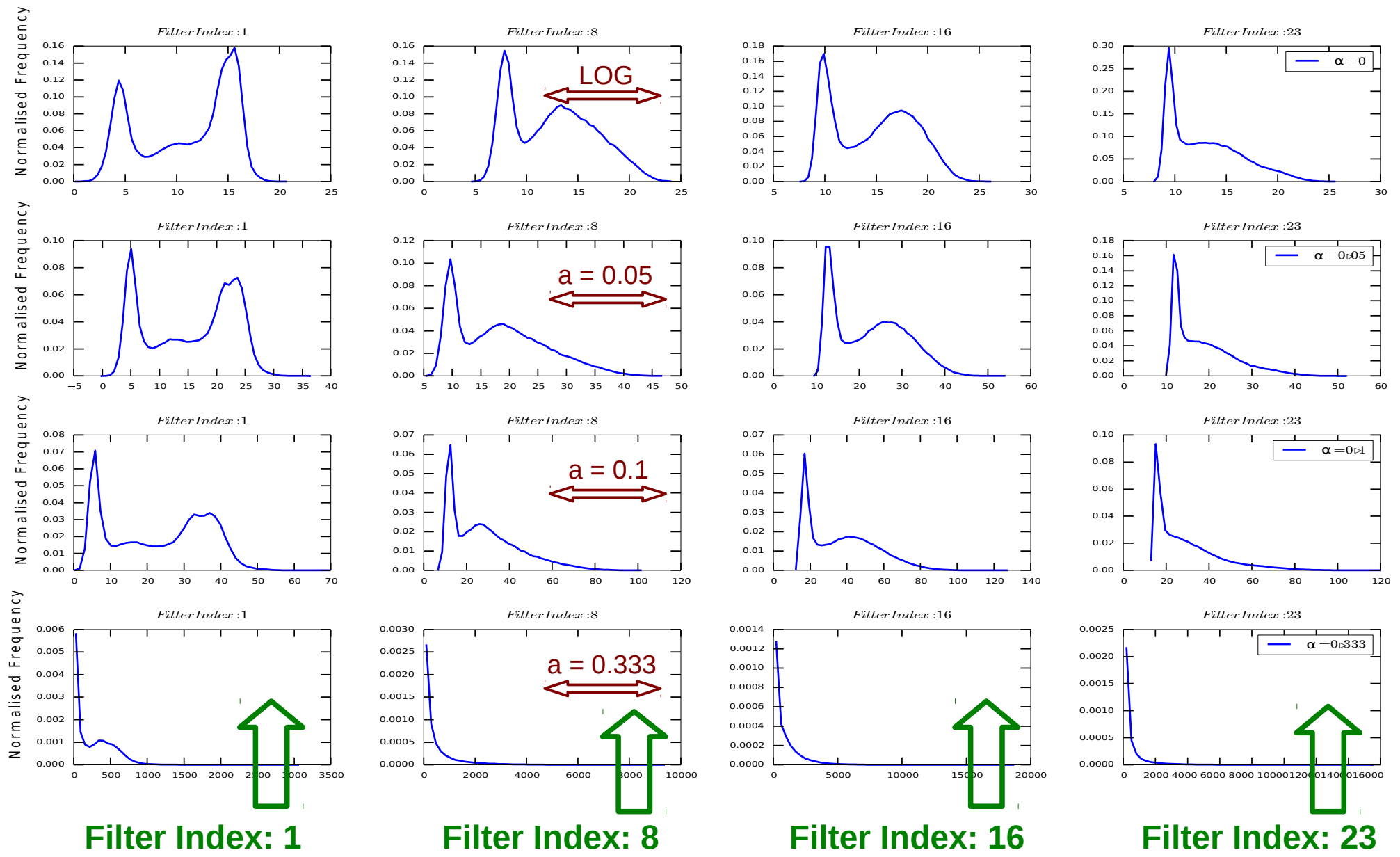
- 1) Statistical effect of GenLog
- 2) GenLog can improve the SNR
- 3) Parameter adjustment in gVTS framework
- 4) Channel noise estimation
 - 1) Pseudocode
 - 2) Initialisation and iteration effects
- 5) DNN Architecture → Bottleneck
- 6) Super-additivity vs Sub-additivity





- NBins: 50
- 330 Utterances, WSJ
- #frames > 241 k

Statistical Effect of GenLog



Filter Index: 1

Filter Index: 8

Filter Index: 16

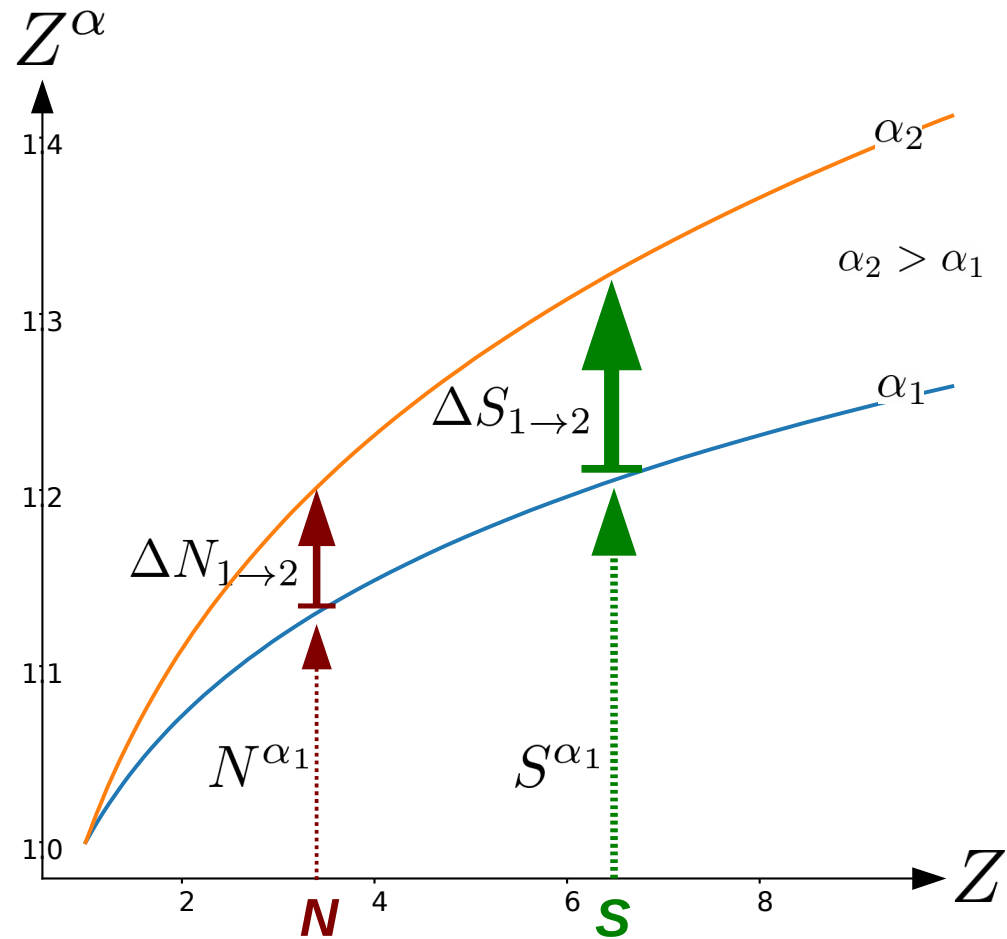
Filter Index: 23



GenLog can improve the SNR ...

$$SNR_1 = \frac{S^{\alpha_1}}{N^{\alpha_1}}$$

$$SNR_2 = \frac{S^{\alpha_2}}{N^{\alpha_2}} = \frac{S^{\alpha_1} + \Delta S_{1 \rightarrow 2}}{N^{\alpha_1} + \Delta N_{1 \rightarrow 2}}$$



$$\alpha_1 < \alpha_2 \Rightarrow SNR_1 < SNR_2$$



$$\begin{cases} \text{GenLog}(x; \alpha) = \frac{1}{\alpha}(x^\alpha - 1), & x > 0 \quad \alpha \neq 0 \\ \lim_{\alpha \rightarrow 0} \text{GenLog}(x; \alpha) = \log(x) \end{cases}$$

↓ α

**Statistical
Distribution**

Easier to be fitted
by a GMM!



↑ α

**SNR
Boost**

$$0.05 \leq \alpha < 0.1$$





Channel Estimation Pseudocode

0. Initialise \mathbf{H}

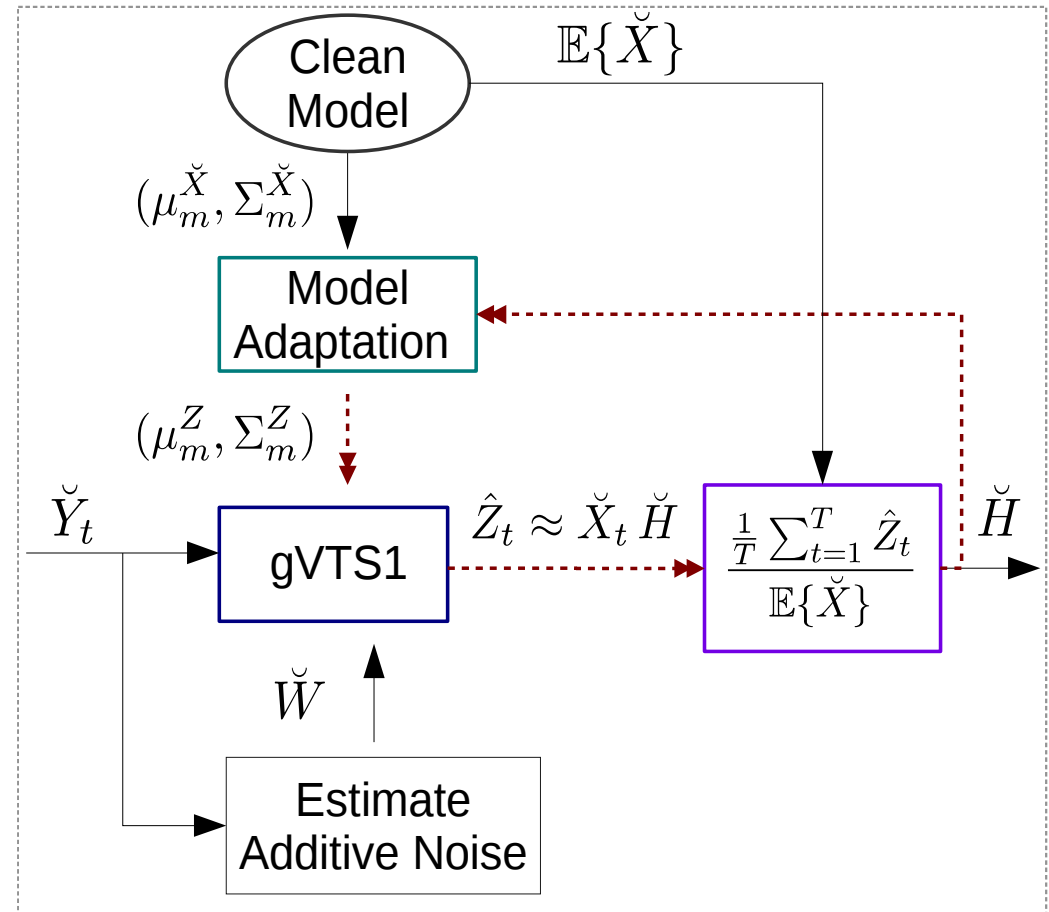
1. Adapt Clean Model with \mathbf{H}

2. gVTS for Additive Noise

3. Update \mathbf{H}

4. If not converged **GO TO 1**

Return \mathbf{H}

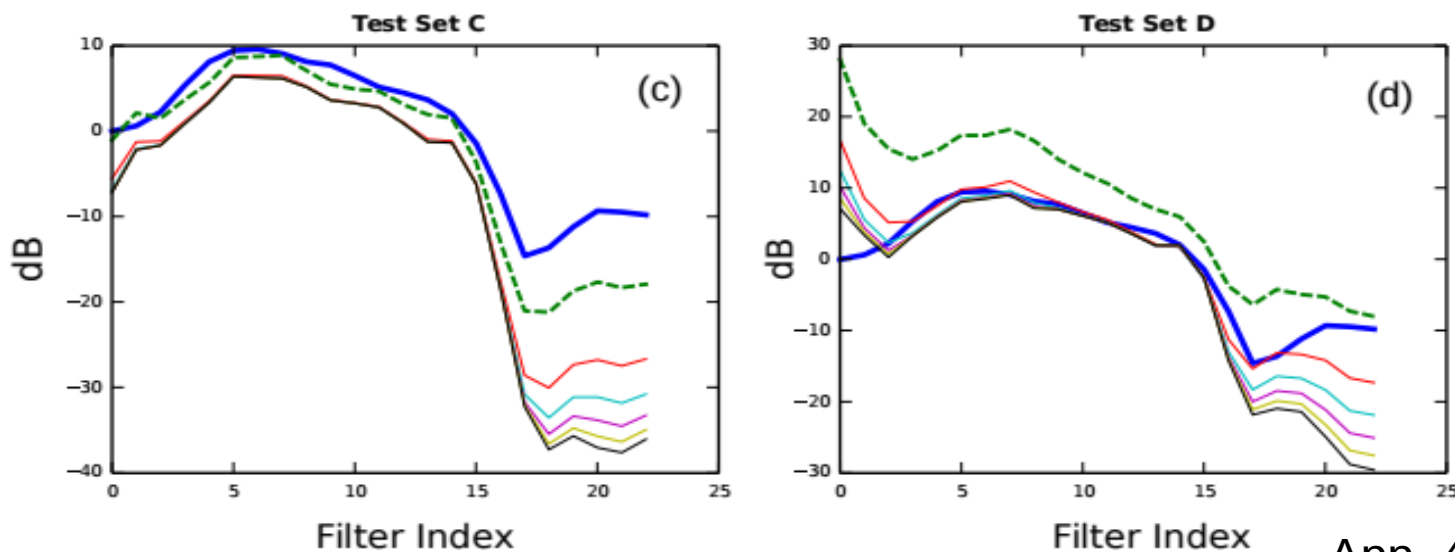
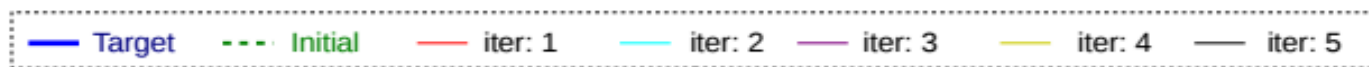
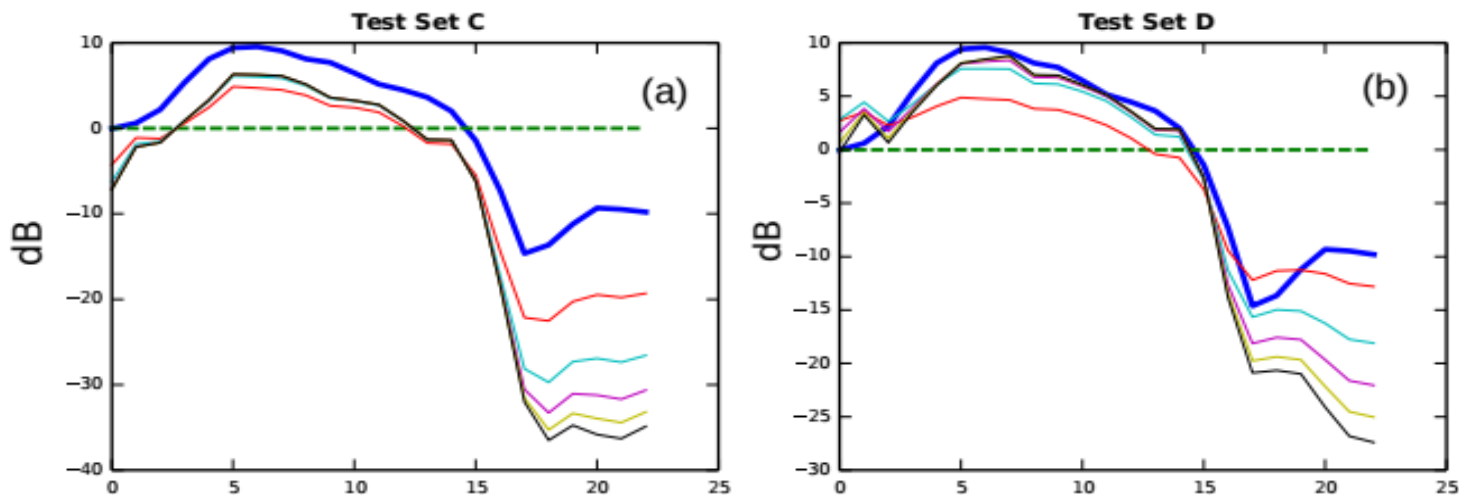




Channel Estimation -- Initialisation and Iteration effects

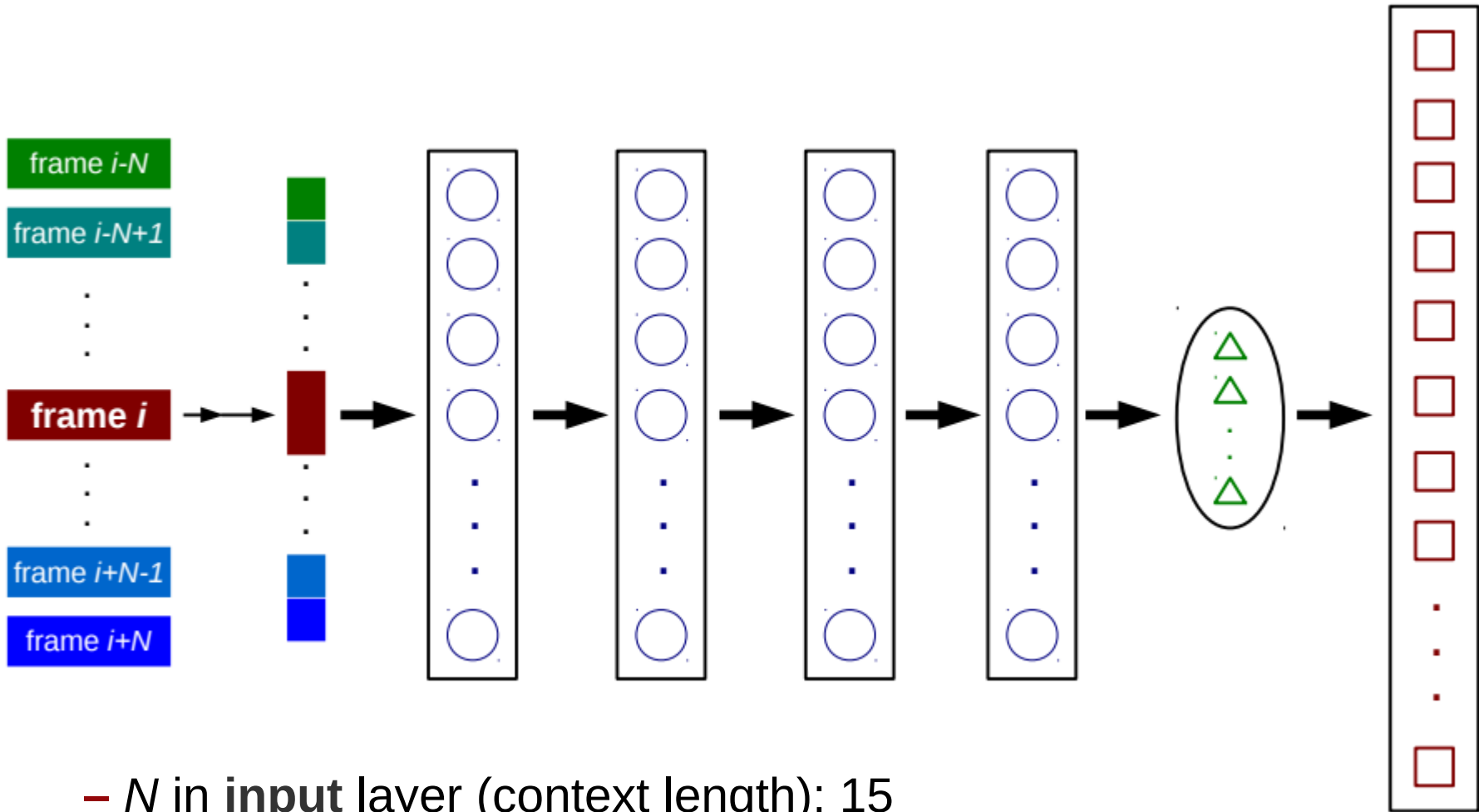
$$H_0 = 1$$

$$H_0 = \frac{Y}{\mathbb{E}\{X\}}$$





DNN Architecture → Bottleneck

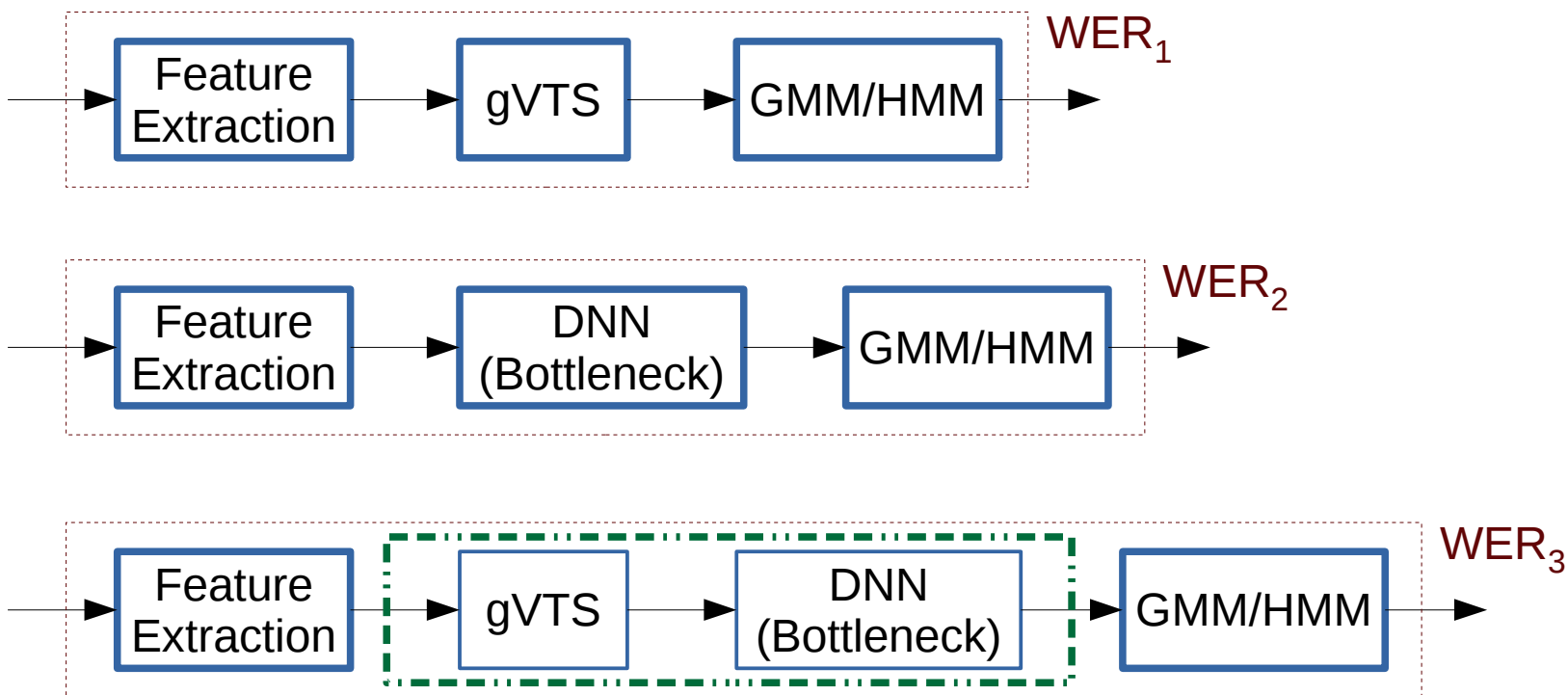


- N in input layer (context length): 15
- #nodes in the hidden layers: 1300
- #nodes in the bottleneck layer: 26
- #nodes in the output layer: state-clustered triphones (~2000)





Super-additivity vs Sub-additivity of a Tandem gVTS-DNN System



- **Super-additive:** WER₃ is better than $\min(\text{WER}_1, \text{WER}_2)$
- **Sub-additive:** WER₃ is worse than $\min(\text{WER}_1, \text{WER}_2)$

