

# Introduction

- Research Group:
  - Jing Hu (Intuit)
  - Xing Han, Joydeep Ghosh (University of Texas at Austin)
- Publications on Hierarchical Time Series Forecasting
  - Simultaneously Reconciled Quantile Forecasting of Hierarchically Related Time Series
  - Dynamic Combination of Heterogeneous Models for Hierarchical Time Series
  - Efficient Forecasting of Large-Scale Hierarchical Time Series via Multilevel Clustering

# Simultaneously Reconciled Quantile Forecasting of Hierarchically Related Time Series (SHARQ)

Han, X., Dasgupta, S., Ghosh, J.

Simultaneously Reconciled Quantile Forecasting of Hierarchically Related Time Series  
International Conference on Artificial Intelligence and Statistics (AISTATS), 2021

# One-Slide Summary

- Motivation:
  - Provide coherent **point and probabilistic multi-step** forecasts for hierarchical time-series
- Approach - SHARQ<sup>1</sup>:
  - Learns **multiple probabilistic forecasts** at the same time
  - **Reconciles probabilistic forecasts simultaneously** during model training
- Empirical Impact:
  - Produced **accurate and coherent** point forecasts
  - Provided **prediction intervals at any specified level**

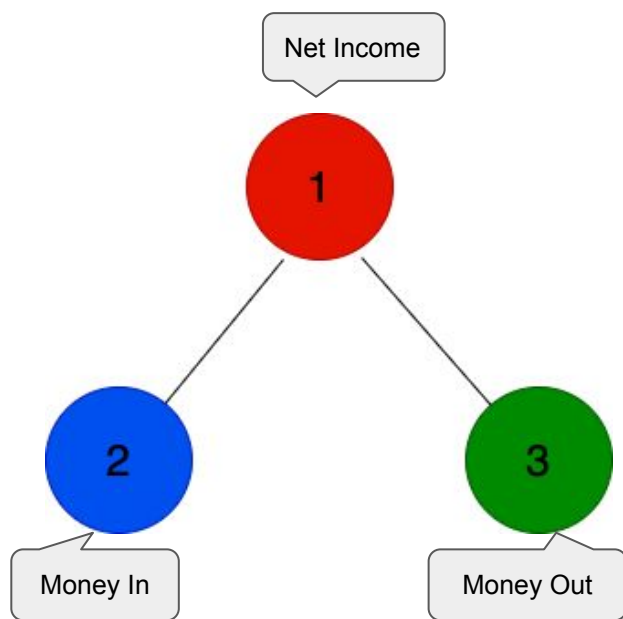
1. SHARQ stands for Simultaneously Hierarchical Reconciled Quantile

# Existing Works: Post Reconciliation



- Post Reconciliation can be posed as OLS
- The reconciled Forecasts are computed as:  $\tilde{\mathbf{y}}_T(h) = \mathbf{S}(\mathbf{S}'\mathbf{\Sigma}_h^+\mathbf{S})^{-1}\mathbf{S}'\mathbf{\Sigma}_h^+\hat{\mathbf{y}}_T(h),$
- Challenges:
  - Only works on assumptions of unbiased base forecasts and Gaussian noise
  - Matrix inversion is not stable and time consuming
  - Forecasts are “forced” to sum up (a hard constraint)
  - Cannot provide prediction intervals

# SHARQ<sup>1</sup>: Reconciling Point Forecast



Estimate Multiple  
Quantiles

Reconciliation

$$L(\hat{y}_1, y_1) = \sum_{\tau = \tau_0}^{\tau_q} \rho_{\tau}(\hat{y}_1, y_1) + \lambda (\hat{y}_1^{50} - \sum_{i \in \{2,3\}} \hat{y}_i^{50})^2$$



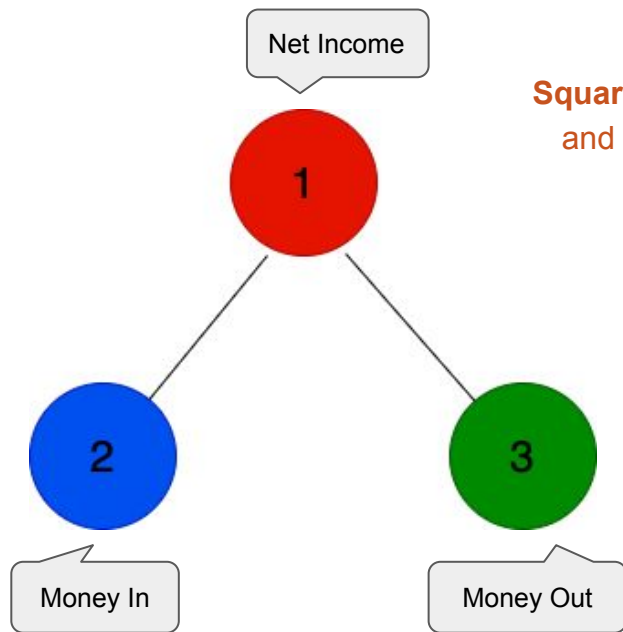
**Bottom Up**

$$L(\hat{y}_i, y_i) = \sum_{\tau = \tau_0}^{\tau_q} \rho_{\tau}(\hat{y}_i^{\tau}, y_i) \quad i = 2, 3$$

Where  $\hat{y}_{t+1}^{\tau} \in \underset{q}{\operatorname{argmin}} \rho_{\tau}(q, y_{t+1}) = \begin{cases} (y_{t+1} - q) \tau & y_{t+1} \geq q \\ (y_{t+1} - q) (\tau - 1) & y_{t+1} < q \end{cases}$

1. SHARQ stands for Simultaneously HierArchical Reconciled Quantile

# SHARQ<sup>1</sup>: Reconciling Quantile Forecasts



$$L(\hat{y}_1, \hat{y}_2, \hat{y}_3) = \sum_{\tau = \tau_0}^{\tau_q} [(\hat{y}_1^\tau - \hat{y}_1^{50})^2 - (\hat{y}_2^\tau - \hat{y}_2^{50})^2 - (\hat{y}_3^\tau - \hat{y}_3^{50})^2]^2$$

Squared distance between each quantiles and median follows additive property in Gaussian distribution.



$$L(\hat{y}_1, y_1) = \sum_{\tau = \tau_0}^{\tau_q} \rho_\tau(\hat{y}_1^\tau, y_1) + \lambda (\hat{y}_1^{50} - \sum_{i \in \{2,3\}} \hat{y}_i^{50})^2$$



$$L(\hat{y}_i, y_i) = \sum_{\tau = \tau_0}^{\tau_q} \rho_\tau(\hat{y}_i^\tau, y_i) \quad i = 2, 3$$

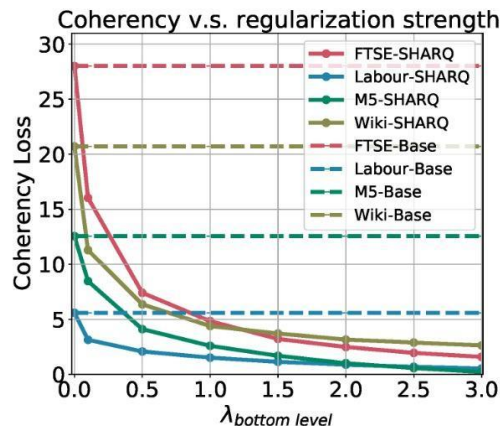
1. SHARQ stands for Simultaneously HierArchical Reconciled Quantile

# Results & Discussions

Algorithm	RNN			
Reconciliation	Level			
	1	2	3	4
BU	15.23	15.88	19.41	<b>17.96</b>
Base	12.89	14.26	16.96	<b>17.96</b>
MinT-sam	14.98	15.94	17.79	19.23
MinT-shr	14.46	15.43	16.94	18.75
MinT-ols	15.01	15.96	18.75	19.21
ERM	14.73	16.62	19.51	20.13
SHARQ	<b>12.55</b>	<b>13.21</b>	<b>16.01</b>	<b>17.96</b>
BU	11.42	12.04	12.32	<b>11.72</b>
Base	10.63	10.15	11.23	<b>11.72</b>
MinT-sam	11.25	11.67	11.87	12.99
MinT-shr	10.76	11.03	11.49	12.81
MinT-ols	11.75	11.56	12.06	13.05
ERM	11.86	12.01	12.42	13.54
SHARQ	<b>9.87</b>	<b>9.68</b>	<b>10.41</b>	<b>11.72</b>

- Explanation of the table:
  - Australian Labour and M5 competition data
  - Performance measured by MAPE
    - The lower, the better
  - Level1 is the top level and level 4 is the bottom level
- Results:
  - The bottom-up approach of Sharq leads to the same performance as BU and Base at the bottom level
  - Performed better than other baseline methods, particularly at higher aggregation levels

# Results & Discussions



Coherency loss drops dramatically after incorporating the regularization

- SHARQ provides a learnable trade-off between coherency and accuracy

Time (s)	FTSE		Labour		M5		Wikipedia	
	training	inference	training	inference	training	inference	training	inference
Base	115.96	0.01	68.35	0.00	181.58	0.00	205.47	0.01
BU	65.83	0.03	57.06	0.00	105.45	0.00	142.53	0.01
MinT-sam	106.55	1,784.77	72.24	430.42	172.11	1,461.81	208.26	1,106.70
MinT-shr	104.35	1,148.49	60.83	317.02	175.83	1,039.53	198.16	788.31
MinT-ols	103.23	1,129.45	64.14	310.13	163.24	977.88	196.88	702.02
ERM	547.66	0.05	497.88	0.01	551.60	0.01	1,299.30	0.04
<b>SHARQ</b>	121.84	0.01	99.96	0.00	201.40	0.00	241.97	0.01

Comparing the average training and inference time across forecasting models

- SHARQ reduces inference time from post-training reconciliation methods



# Dynamic Combination of Heterogeneous Models for Hierarchical Time Series (DYCHEM)

**Hu, J**, Han, X., Ghosh, J.

Dynamic Combination of Heterogeneous Models for Hierarchical Time Series  
International Conference on Data Mining (ICDM), 2022

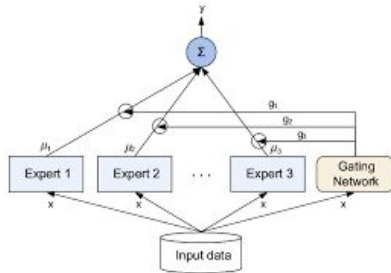
# One-Slide Summary

- Motivations:
  - Borrow strength from heterogeneous forecasting models
  - Allow any time series time algorithm to be used as a base forecasting method
  - Eliminate quantile crossing behavior of quantile regression
- DYCHEM<sup>1</sup>:
  - Improves point forecasts by utilizing multiple forecasting models via a **mixture-of-expert** framework
  - Allows mixing of **different time series forecasting algorithms**
  - Quantile forecasts are coherent and **without quantile crossing**
- Empirical Impact:
  - Significantly improved point & probabilistic forecasting performance

1. DYCHEM stands for DYnamic Combination of Heterogeneous Models

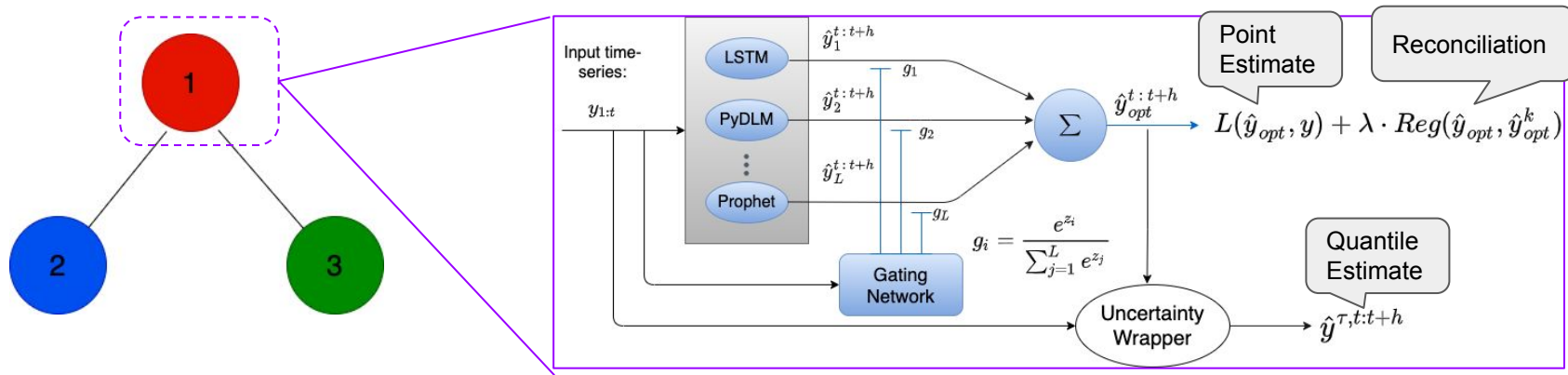
# Motivations

- We wish to borrow the strength of heterogeneous forecasting models (DLM, Prophet, Auto-ARIMA, Deep learning methods), given their expertise on different types of data
- SHARQ requires to modify the objective function of forecasting models, which is not always available for existing state-of-the-art forecasting libraries
- We wish to build a framework where the forecasting models can be independent and user-specified



Mixture-of-Experts: is an **ensemble learning method** that utilize **gating network** to combine outputs from each expert.

# DYCHEM - DYnamic Combination of Heterogeneous Models



Improved SHARQ via:

- Combine heterogeneous models at each vertex.
- Each expert is independently configurable and replaceable; overall is a black-box model.
- How about quantile estimations? Simply combining distributions won't work.

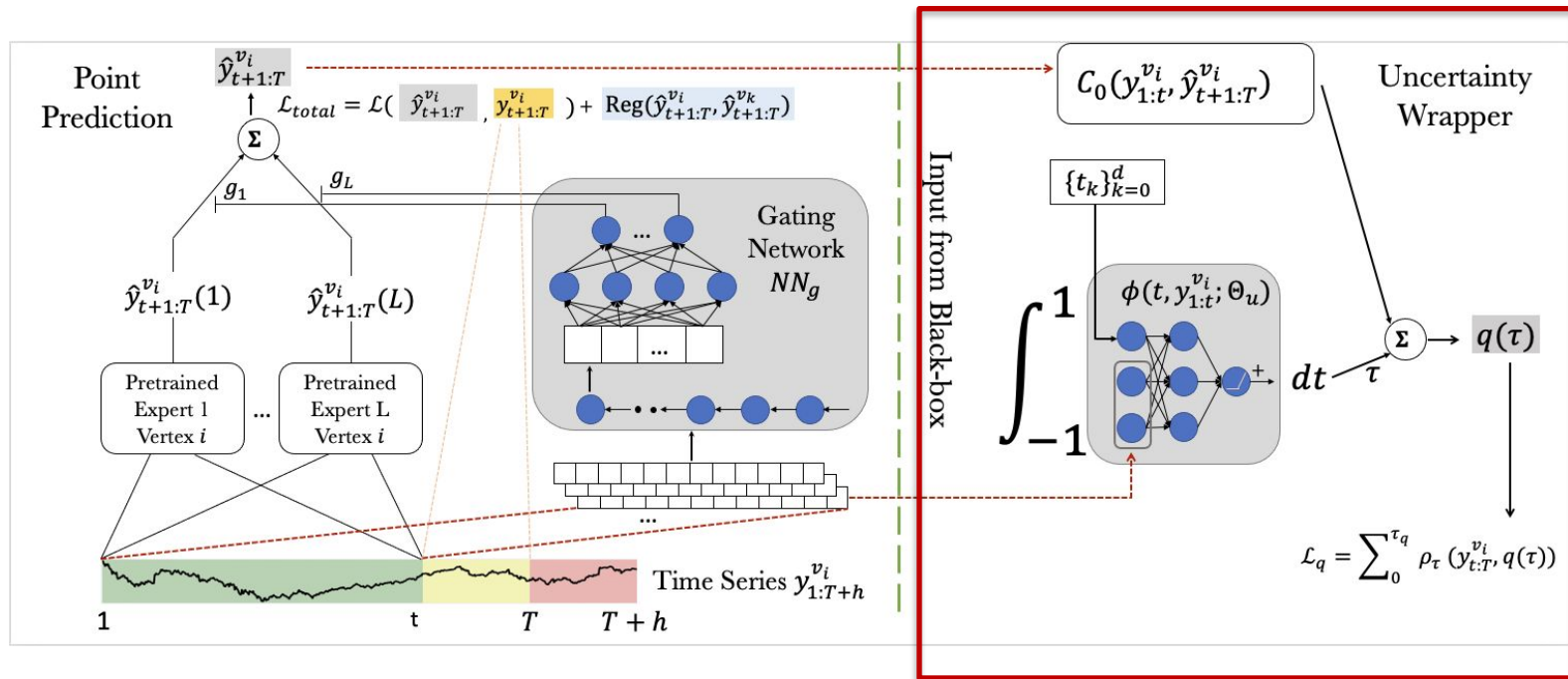
# Overall Structure

Requires 3 inputs:

- Point prediction
- Time series training data
- Quantile levels

Outputs:

- Quantile estimations at specified levels



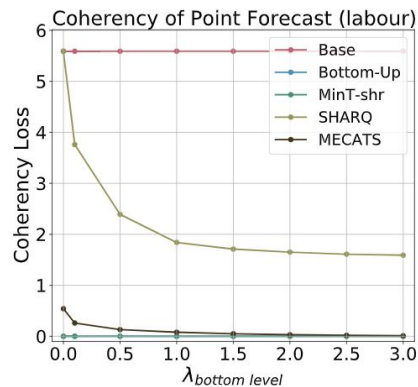
- The heterogeneous experts bring robustness to the prediction.

# DYCHEM: Comparison with Baselines

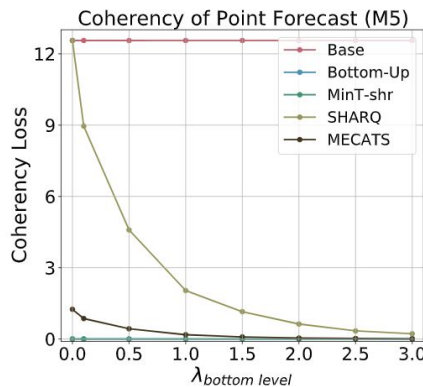
Data \ Method	Level	DYCHEM-LOO	SHARQ	HIER-E2E	Average	DYCHEM
Labour	1	43.33 $\pm$ 0.42 (.054)	52.07 $\pm$ 0.45 (.085)	45.12 $\pm$ 0.23 (.085)	49.34 $\pm$ 0.65 (.075)	<b>38.84 <math>\pm</math>0.04 (.045)</b>
	2	53.68 $\pm$ 0.68 (.104)	58.69 $\pm$ 0.41 (.120)	55.61 $\pm$ 0.74 (.107)	60.87 $\pm$ 0.33 (.119)	<b>48.64 <math>\pm</math>0.78 (.092)</b>
	3	57.16 $\pm$ 0.25 (.135)	64.02 $\pm$ 0.09 (.132)	60.03 $\pm$ 0.26 (.134)	69.29 $\pm$ 0.42 (.138)	<b>49.17 <math>\pm</math>0.36 (.144)</b>
	4	65.05 $\pm$ 0.18 (.153)	72.13 $\pm$ 0.34 (.167)	71.38 $\pm$ 0.15 (.154)	75.56 $\pm$ 0.94 (.156)	<b>61.22 <math>\pm</math>0.14 (.163)</b>
M5	1	49.29 $\pm$ 0.34 (.071)	56.31 $\pm$ 0.17 (.054)	51.69 $\pm$ 0.05 (.070)	59.61 $\pm$ 0.38 (.104)	<b>42.61 <math>\pm</math>0.14 (.046)</b>
	2	54.36 $\pm$ 0.28 (.127)	62.16 $\pm$ 0.27 (.079)	54.72 $\pm$ 0.63 (.116)	60.48 $\pm$ 0.58 (.133)	<b>49.75 <math>\pm</math>0.22 (.084)</b>
	3	55.18 $\pm$ 0.22 (.142)	65.37 $\pm$ 0.63 (.134)	65.02 $\pm$ 0.24 (.142)	68.29 $\pm$ 0.25 (.143)	<b>53.61 <math>\pm</math>0.42 (.101)</b>
	4	59.04 $\pm$ 0.36 (.164)	72.86 $\pm$ 0.27 (.189)	72.04 $\pm$ 0.36 (.164)	70.29 $\pm$ 0.34 (.168)	<b>57.89 <math>\pm</math>0.47 (.109)</b>
AEDemand	1	61.35 $\pm$ 0.76 (.132)	68.19 $\pm$ 0.29 (.113)	64.45 $\pm$ 0.48 (.213)	67.32 $\pm$ 0.29 (.164)	<b>59.89 <math>\pm</math>0.32 (.145)</b>
	2	58.12 $\pm$ 0.46 (.152)	66.57 $\pm$ 0.24 (.199)	63.72 $\pm$ 0.36 (.131)	63.58 $\pm$ 0.72 (.129)	<b>55.72 <math>\pm</math>0.73 (.122)</b>
	3	66.38 $\pm$ 0.78 (.124)	68.25 $\pm$ 0.47 (.131)	68.01 $\pm$ 0.22 (.126)	70.44 $\pm$ 0.09 (.124)	<b>62.55 <math>\pm</math>0.14 (.111)</b>
	4	76.58 $\pm$ 0.63 (.136)	87.35 $\pm$ 0.69 (.225)	82.47 $\pm$ 0.28 (.192)	73.22 $\pm$ 0.37 (.135)	<b>71.45 <math>\pm</math>0.43 (.125)</b>
Wiki	1	65.98 $\pm$ 0.22 (.121)	70.36 $\pm$ 0.24 (.147)	69.67 $\pm$ 0.58 (.067)	66.42 $\pm$ 0.16 (.128)	<b>63.27 <math>\pm</math>0.73 (.117)</b>
	2	68.54 $\pm$ 0.47 (.157)	73.06 $\pm$ 0.42 (.159)	68.24 $\pm$ 0.33 (.108)	72.01 $\pm$ 0.52 (.157)	<b>65.14 <math>\pm</math>0.46 (.143)</b>
	3	72.42 $\pm$ 0.36 (.149)	76.15 $\pm$ 0.34 (.135)	74.62 $\pm$ 0.19 (.155)	74.37 $\pm$ 0.83 (.147)	<b>69.48 <math>\pm</math>0.33 (.156)</b>
	4	77.12 $\pm$ 0.23 (.268)	78.42 $\pm$ 0.34 (.201)	79.63 $\pm$ 0.41 (.291)	81.38 $\pm$ 0.65 (.278)	<b>75.69 <math>\pm</math>0.76 (.189)</b>
	5	84.77 $\pm$ 0.49 (.241)	85.12 $\pm$ 0.62 (.345)	79.65 $\pm$ 0.24 (.326)	84.68 $\pm$ 0.42 (.221)	<b>76.88 <math>\pm</math>0.72 (.213)</b>

- Forecasting performance measured by averaged MASE and CRPS (within bracket). All experiments are repeated 5 times.
- Baselines: simple averaging (Average); SHARQ; leave-one-out of one expert; HIER-E2E [1]

# DYCHEM: Coherency Analysis

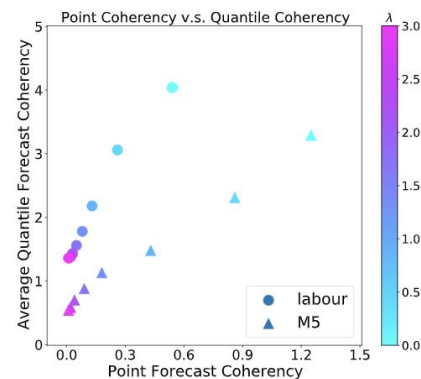


(a)



(b)

$$\sum_{\tau=\tau_0}^{\tau_q} [(\hat{y}_1^\tau - \hat{y}_1^{50})^2 - (\hat{y}_2^\tau - \hat{y}_2^{50})^2 - (\hat{y}_3^\tau - \hat{y}_3^{50})^2]^2$$



(c)

(a), (b) Coherent loss of point forecast w.r.t. regularization strength  $\lambda$  on two datasets. (c) Relationship between point forecast coherency and average of quantile coherency

# Generating Non-Crossing Probabilistic Forecasts

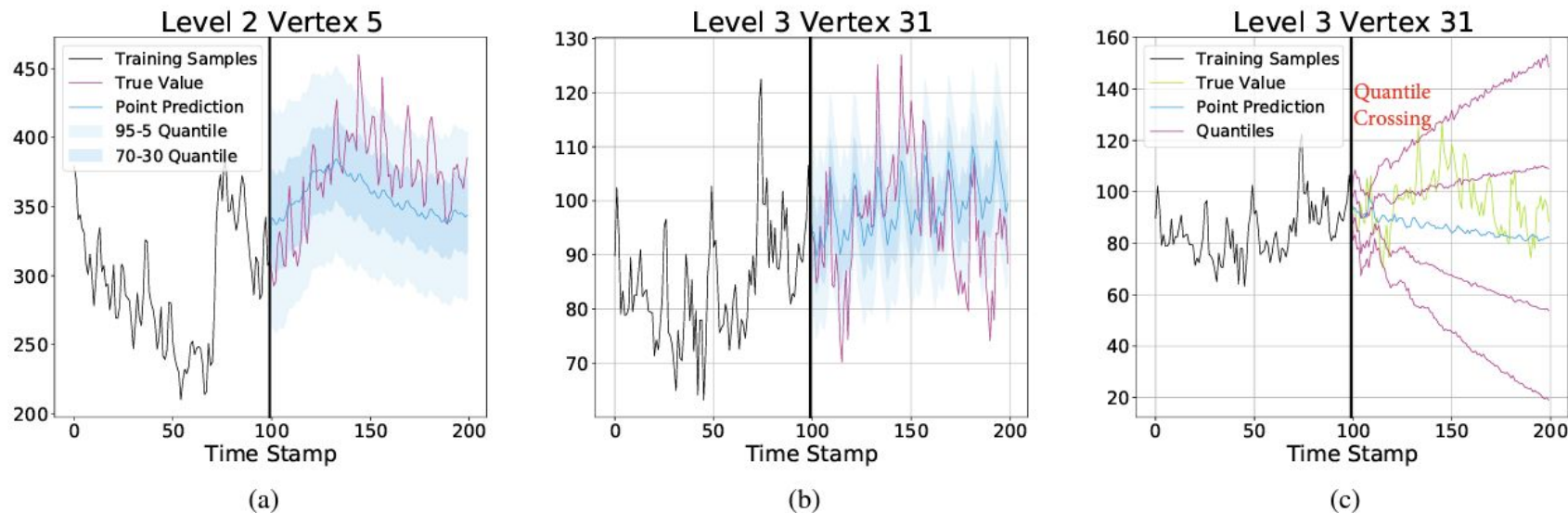


Fig. 4. (a), (b) Quantile forecasting results generated by DYCHEM at vertex 5 and 31 of the Australian Labour data, where  $\tau_s = [0.05, 0.3, 0.5, 0.7, 0.95]$ . (c) SHARQ at same  $\tau_s$ , results showing mild quantile crossing.



# Efficient Forecasting of Large-Scale Hierarchical Time Series via Multilevel Clustering

Hu, J., Han, X., Ren, T., Ghosh, J., Ho, N.

Efficient Forecasting of Large-Scale Hierarchical Time Series via Multilevel Clustering  
International Conference on Time Series and Forecasting, 2023

# One-Slide Summary

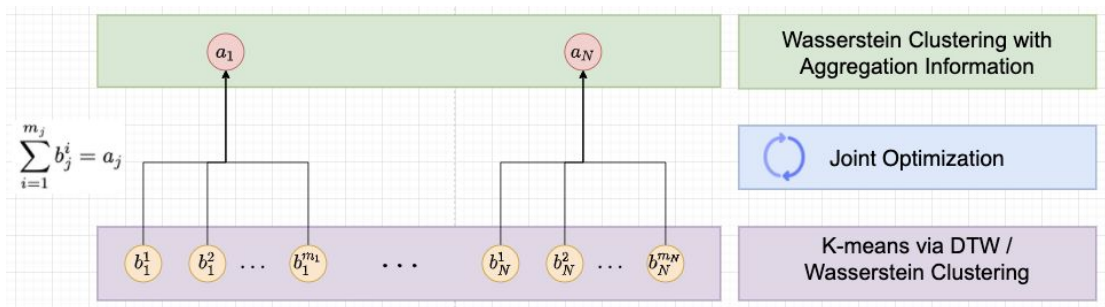
- Motivations:
  - Accelerate large-scale hierarchical time-series (HTS) forecasting
  - Provide analytical insights to large-scale user records with multi-level structure
- Hierarchical Time Series Clustering:
  - Leveraged local aggregation information to improve time series clustering
  - Enables multi-level time-series clustering with arbitrary lengths
- Empirical Values:
  - Provided superior clustering results for time-series with multilevel structures
  - Improved forecasting efficiency for large-scale hierarchical time-series without too much sacrifice on accuracy.

# Existing Approaches

- Mainly focused on time series data without hierarchical structure
  - Distance based approach
    - Define an appropriate distance measurement and cluster based on the distance
    - Extract temporal information capturing features and cluster in the embedding space
  - Model based approach
    - Specify the generative model type and estimate the parameters using maximum likelihood estimation

# High Level View of The Approach

- Our approach to a 2-level hierarchical time-series data:
  - Cluster bottom level data via K-means or Wasserstein Clustering
  - Assign a probability measure to each aggregated-level time series according to the cluster assignment of the child node.
    - Leverage local information in hierarchical time series clustering
  - Cluster top level data via Wasserstein Clustering incorporating the aggregation information
  - Using Wasserstein Clustering, both levels can be simultaneously clustered through joint optimization



# Accelerating Hierarchical Forecasting via Clustering

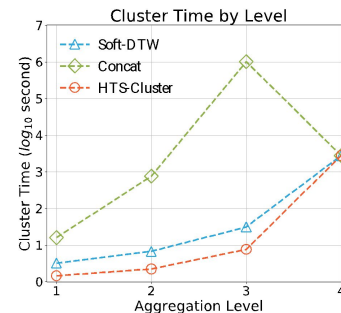
- Forecasts for individual hierarchical time series can “borrow strength” from the forecasts of nearest cluster means at each level.
- Utilizing fuzzy clustering, forecast of each time series can be represented as the weighted combination of forecasts of corresponding cluster means.
- One can quickly finetune results of each time-series to tailor the individual specifics.

# Simulation & E-Commerce Clustering Results

Simulation

Real-world

Method\Metric	Time (s)	Global			Local		
		NMI	AMI	ARI	NMI	AMI	ARI
DTCR	132	0.325 $\pm$ 0.012	0.257 $\pm$ 0.023	0.21 $\pm$ 0.011	0.392 $\pm$ 0.014	0.313 $\pm$ 0.006	0.284 $\pm$ 0.009
Soft-DTW	67	0.412 $\pm$ 0.009	0.326 $\pm$ 0.019	0.277 $\pm$ 0.008	0.411 $\pm$ 0.022	0.342 $\pm$ 0.009	0.304 $\pm$ 0.014
Concat	186	0.436 $\pm$ 0.015	0.342 $\pm$ 0.014	<b>0.314</b> $\pm$ 0.016	0.411 $\pm$ 0.022	0.342 $\pm$ 0.009	0.304 $\pm$ 0.014
<b>HTS-Cluster</b>	37	<b>0.455</b> $\pm$ 0.018	<b>0.354</b> $\pm$ 0.015	0.302 $\pm$ 0.013	<b>0.424</b> $\pm$ 0.018	<b>0.366</b> $\pm$ 0.013	<b>0.321</b> $\pm$ 0.018
DTCR	72	0.065 $\pm$ 0.002	0.015 $\pm$ 0.001	0.008 $\pm$ 0.002	0.105 $\pm$ 0.011	0.059 $\pm$ 0.002	0.054 $\pm$ 0.003
Soft-DTW	49	0.119 $\pm$ 0.005	0.043 $\pm$ 0.003	0.027 $\pm$ 0.003	0.126 $\pm$ 0.008	<b>0.082</b> $\pm$ 0.006	0.061 $\pm$ 0.005
Concat	174	<b>0.135</b> $\pm$ 0.004	0.073 $\pm$ 0.007	<b>0.045</b> $\pm$ 0.006	0.126 $\pm$ 0.008	<b>0.082</b> $\pm$ 0.006	0.061 $\pm$ 0.005
<b>HTS-Cluster</b>	34	0.134 $\pm$ 0.005	<b>0.075</b> $\pm$ 0.005	0.041 $\pm$ 0.004	<b>0.128</b> $\pm$ 0.014	0.064 $\pm$ 0.005	<b>0.065</b> $\pm$ 0.002



- Baseline method
  - DTCR: Deep Temporal Clustering Representation
    - Regular multivariate time series clustering **without considering hierarchical structure.**
- Competing methods
  - Soft-DTW: Soft-DTW divergence-based K-means **on all levels.**
  - Concat: simply **concatenate local time-series** for global time series.
  - HTS-Cluster: Our method, jointly optimizing two levels.

# Accelerating Hierarchical Forecasting

Method \ Level	1	2	3	4	Total Time
Without cluster	62.39	76.26	78.25	84.14	1
DTCR	82.35	96.09	104.85	104.33	0.39
Soft-DTW	78.61	93.04	93.12	96.76	0.27
Concat	74.24	84.65	83.73	96.76	0.57
HTS-Cluster	72.99	80.07	85.29	96.76	0.16

Forecasting massive HTS with the help of clustering, results are measured by mean absolute scaled error (MASE) and relative computing time.

- Jointly model similar HTS via clustering is more efficient than building independent models without sacrificing too much accuracy.

*Thank you!*