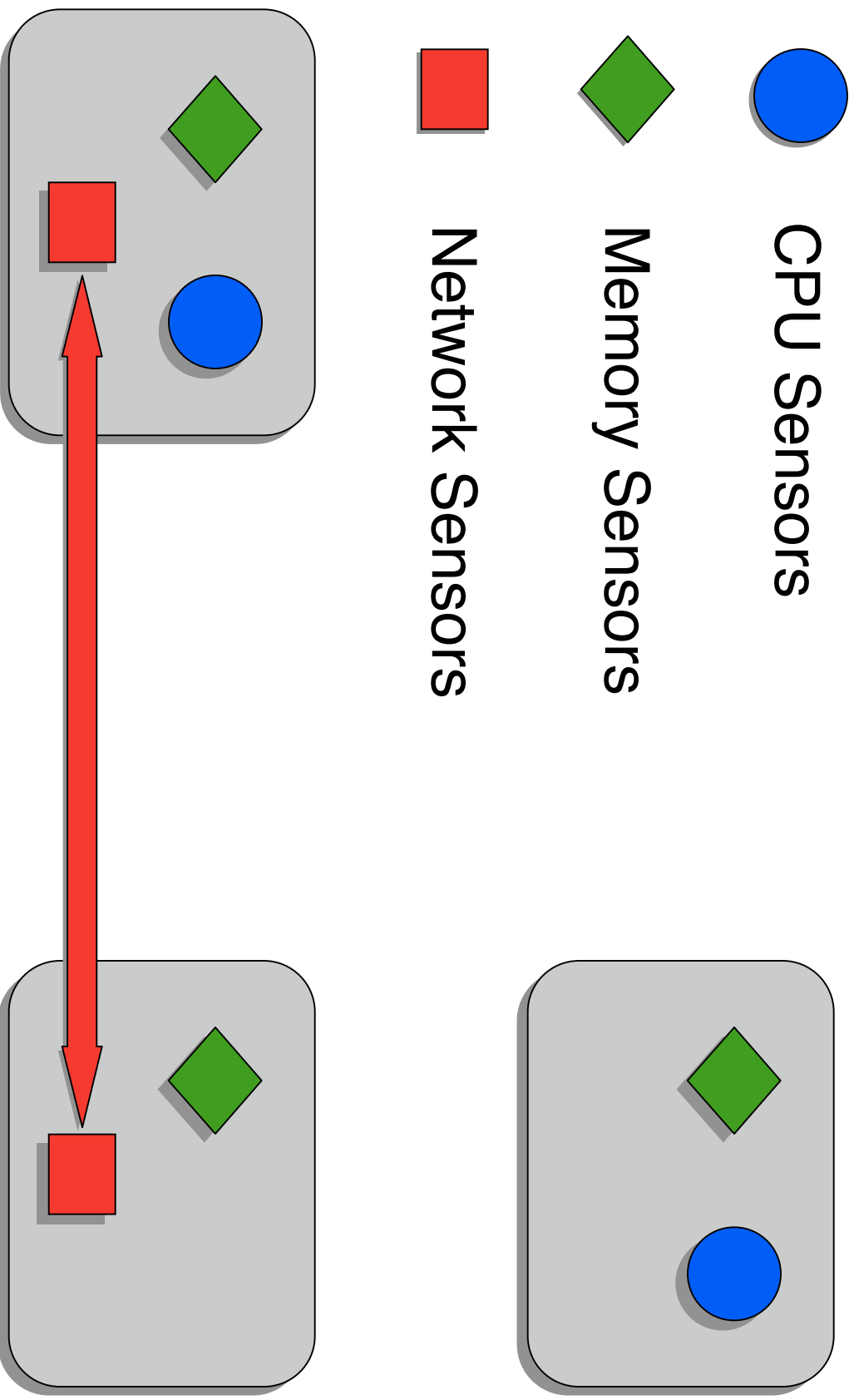


Multivariate Resource Performance Prediction in the NWS

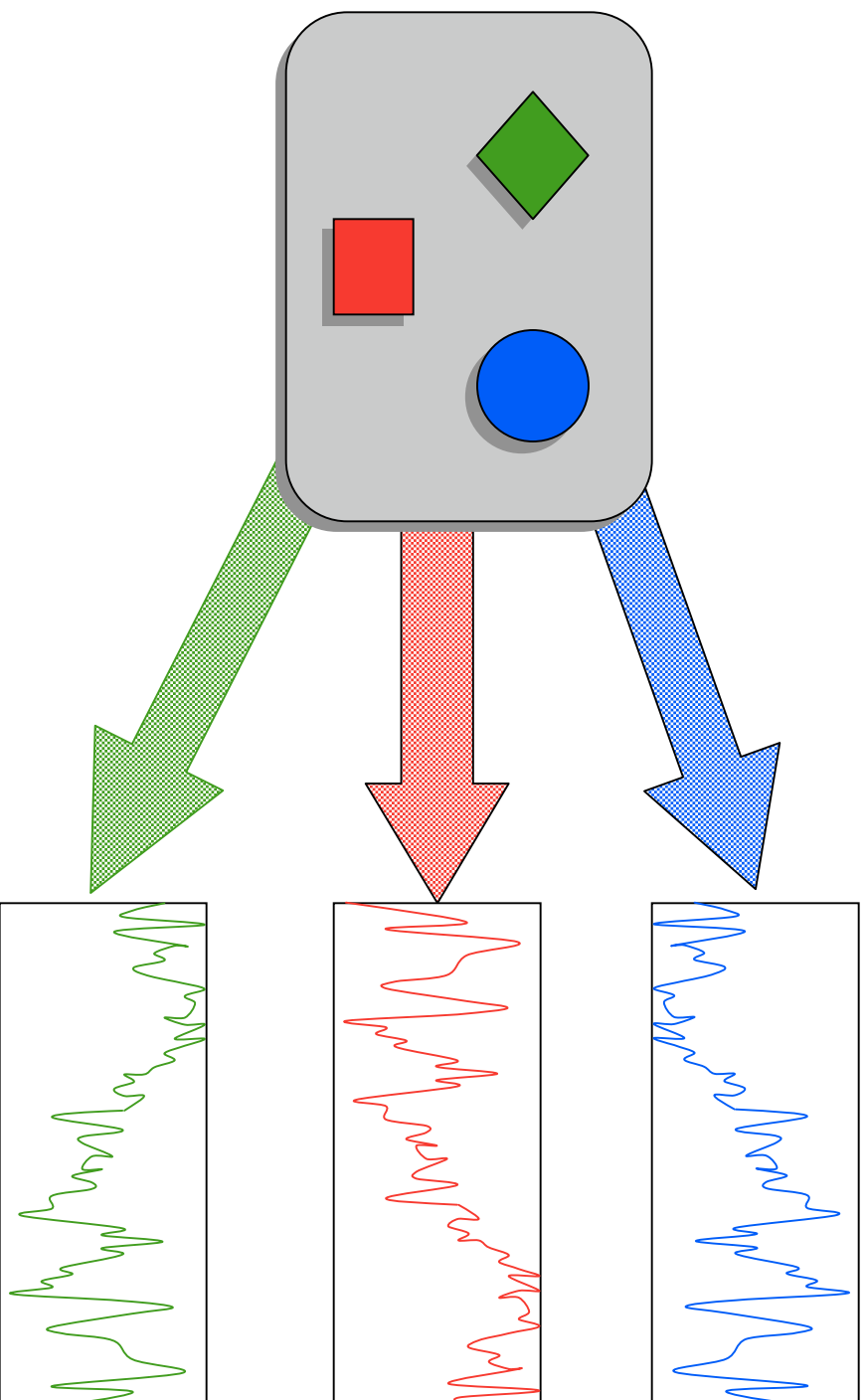
Martin Swamy

UCSB

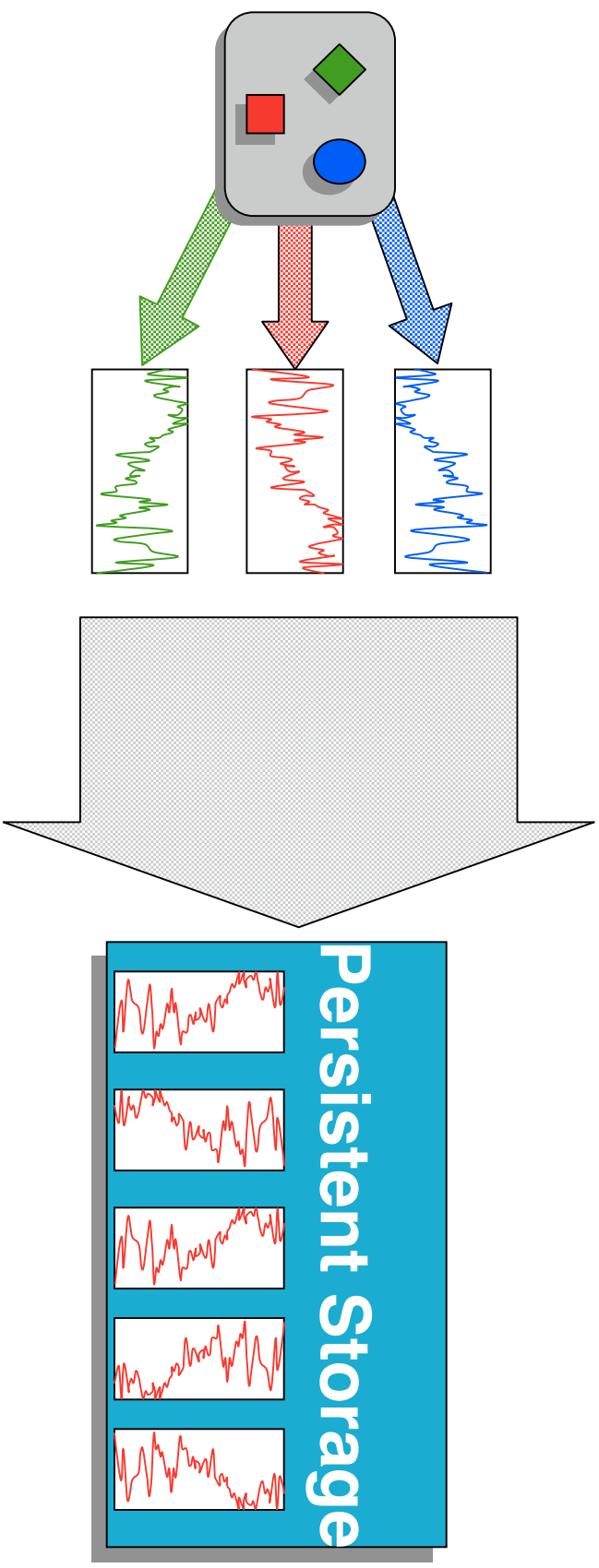
NWS Architecture



Sensor Elements Produce Time Series

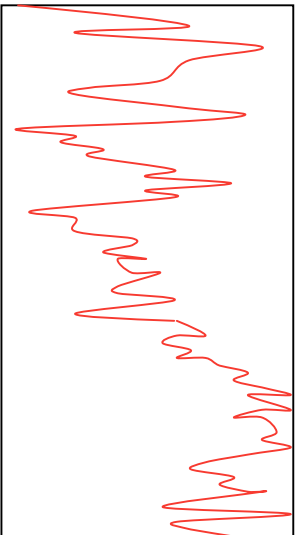


Time Series are Stored

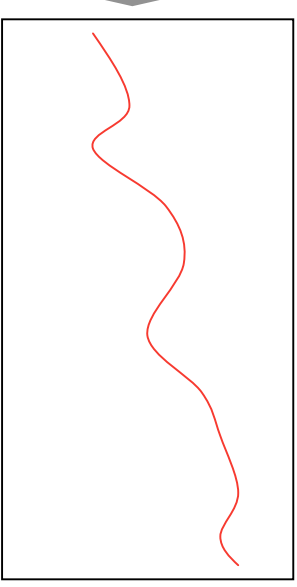


In the `nws_memory`

Time Series are Passed through Forecasters

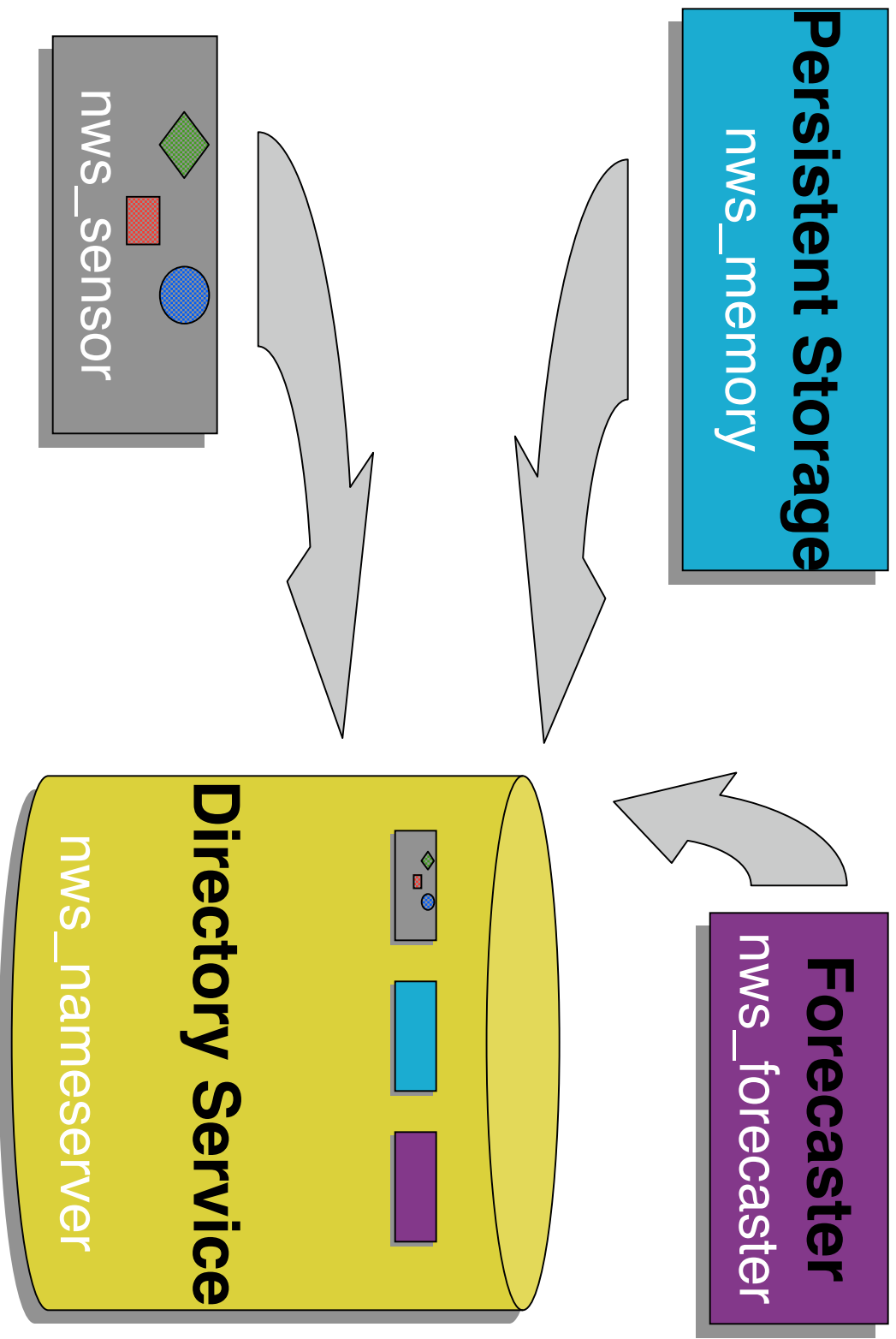


**Forecasting
System**

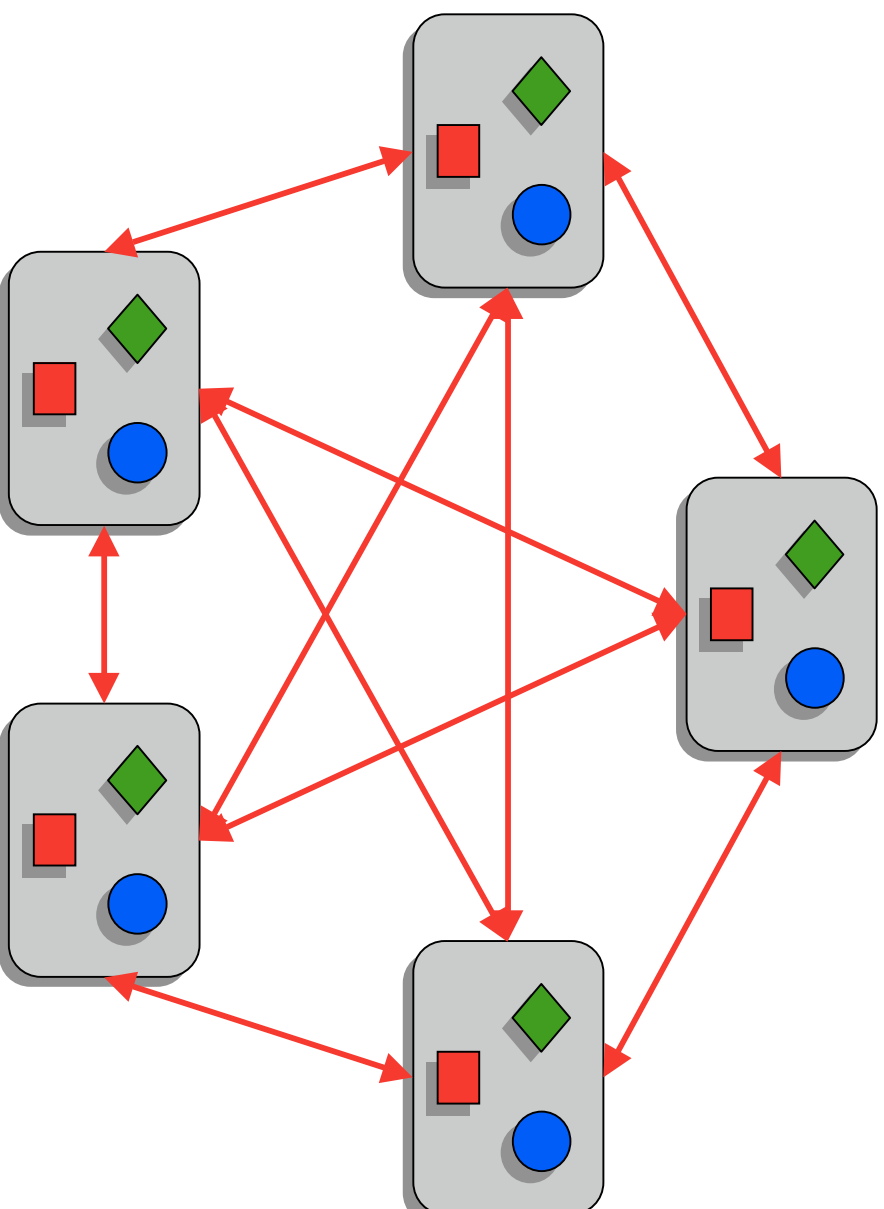


Implemented as a library or a
`daemon - nws_forecaster`

Elements Register with the Directory Service

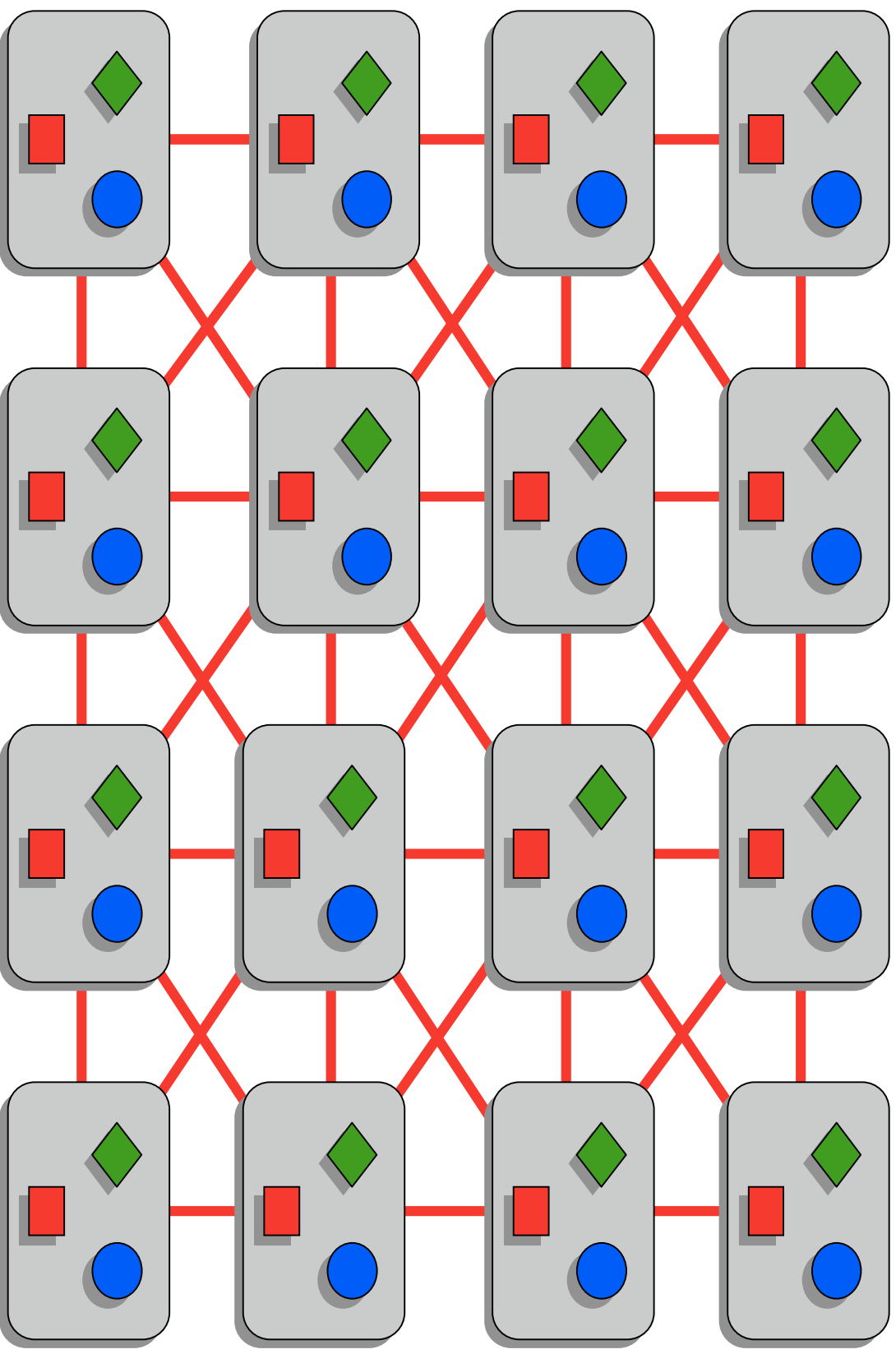


Network Sensors Form Peer Groups

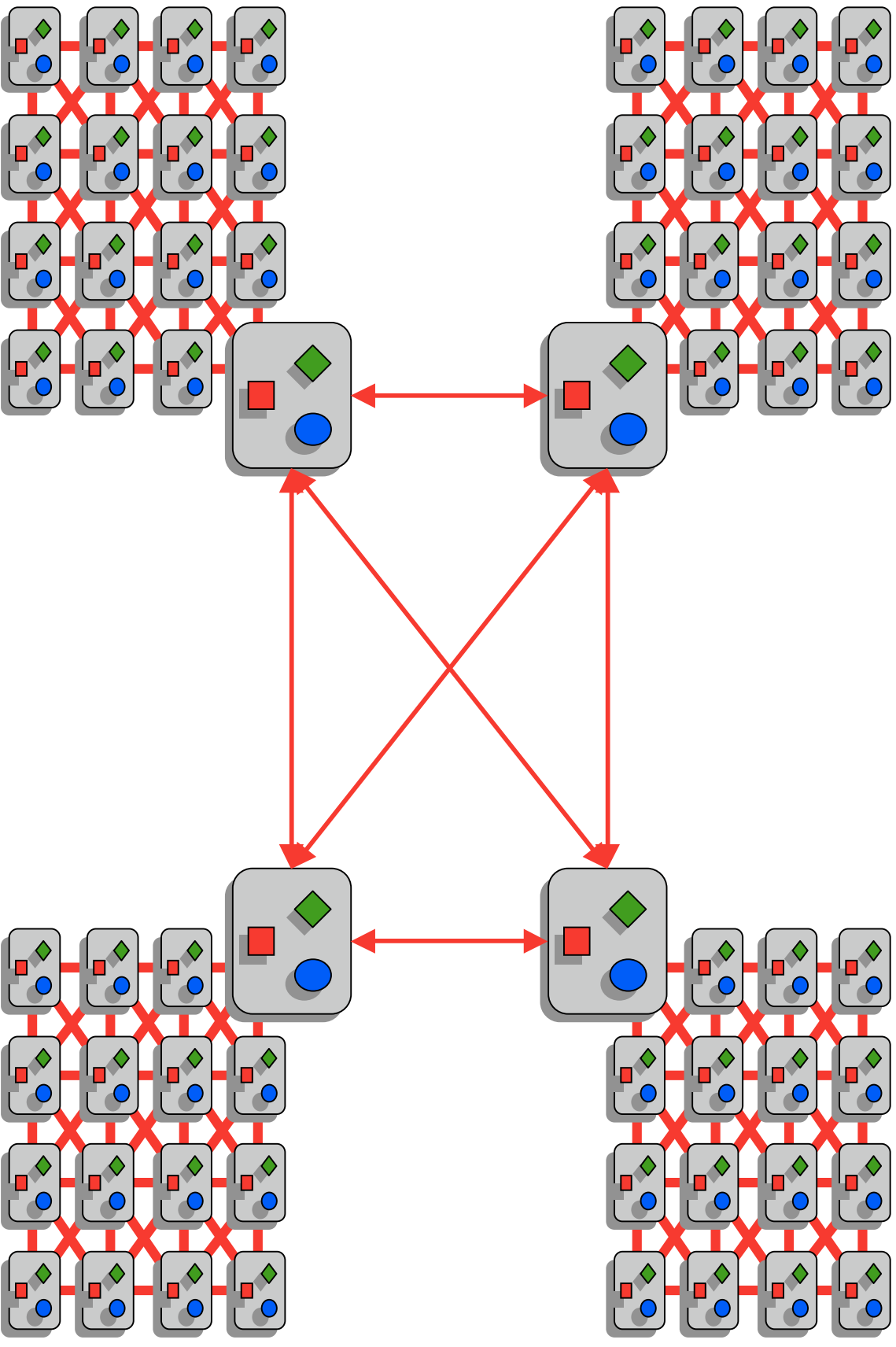


Called “cliques”

A Fully-Connected Cluster or Site



A Hierarchy of Cliques for scalability



Measurement Issues

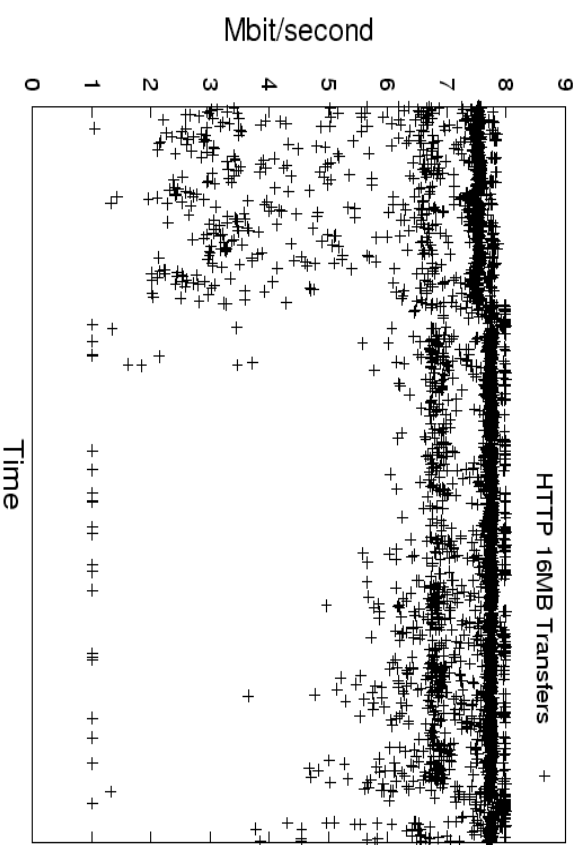
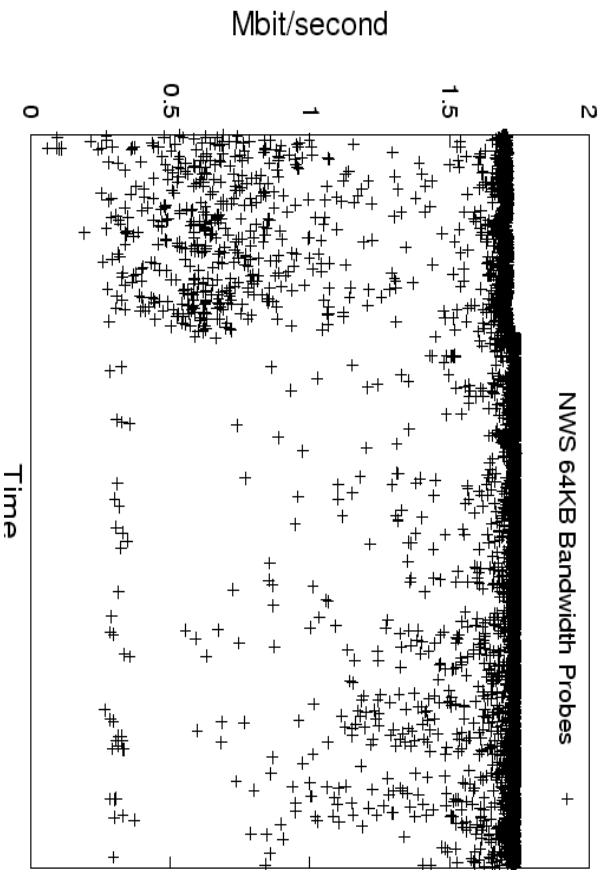
- Active measurements are intrusive
- More accuracy often means more intrusiveness
- The NWS defaults to lightweight probes that have questionable usefulness
 - 64KB transfers
 - Clearly not a measure of nominal bandwidth!
 - Especially as $bw * delay$ grows
- Nominal bandwidth isn't the idea here

Some philosophy

- Resource performance is extremely dynamic
- Presenting up to date performance information to distributed systems is essential for their effective operation
 - Particularly the case for Grid environments
- Most measurement in this context is used as a prediction of future performance
- No privileged access

How much Information?

- Intuitively, there seems to be a relationship between shorter and longer measurements
- Go ahead, squint
- But, they are significantly different

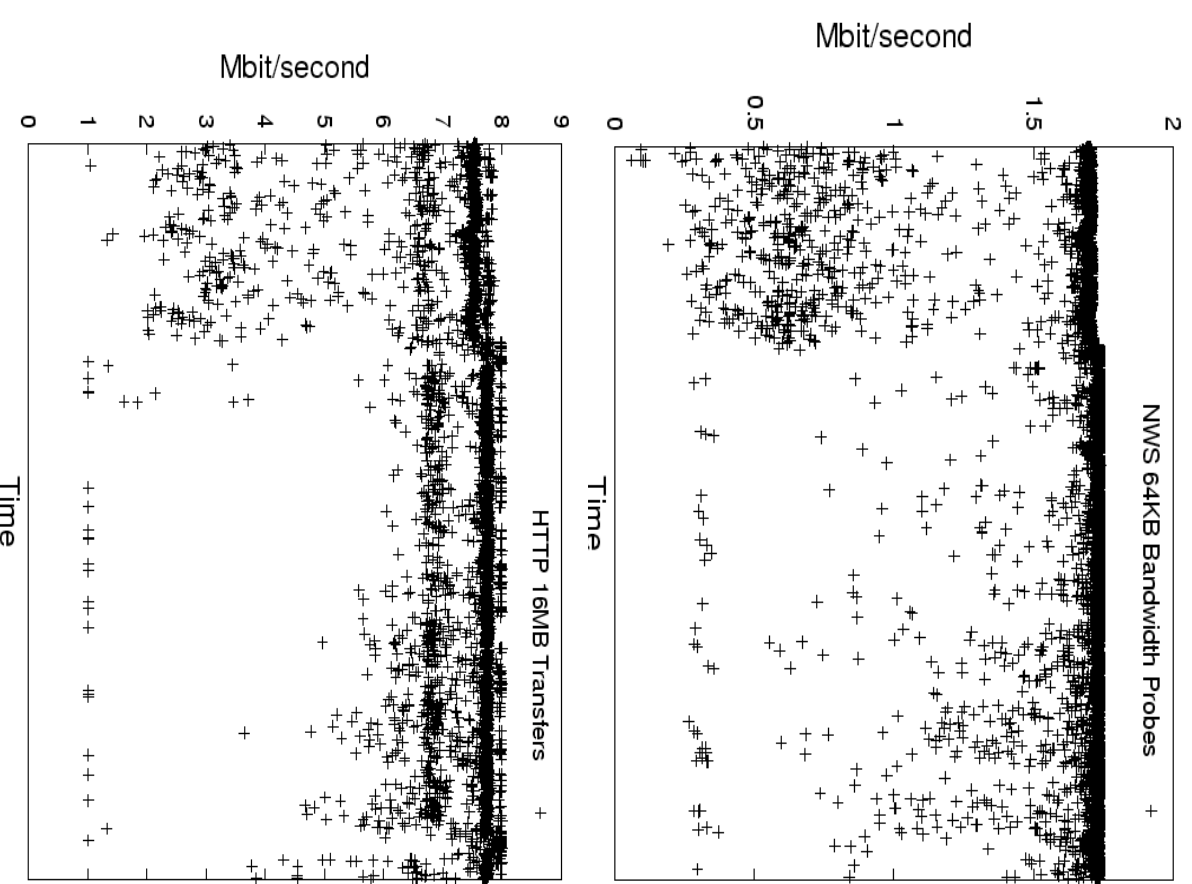


Experimental Methodology

- Collect 64KB NWS measurements every 10 seconds
- Time 16MB HTTP transfers every 60 seconds
 - Using wget
 - The file to be transferred comes from the filesystem
- Heavily used, general purpose system
 - It can't all be in the buffer cache

Recall this picture

- Factor of 4 (ish) difference in maximum throughput
- Figure out the difference and use a simple linear scaling function?



The relationship isn't that easy

- The relationship isn't necessarily linear
 - This makes the regression more difficult
- The relationship might change over time
- The problem of data matching is tricky, too
- Perhaps different amounts of information
 - Particularly if application instrumentation data is used.

Recent Work

- Experiments of Network Throughput Measurement and Forecasting Using the Network Weather Service, P. Primet, R. Harakaly, F. Bonnassieux (INRIA, ENS-Lyon), CCGrid 02
- Attempts to compute the relationship between NWS data and Iperf data with a magic scaling factor
 - We tried this too, but weren't happy with the results

Recent Work

- Predicting Sporadic Grid Data Transfers, S. Vazhkudai, J. Schopf, to appear HPDC 11
- Focuses on using NWS data to predict GridFTP transfers
- Explores a variety of matching techniques for regression

Multivariate Forecasters

- **We want to take a suspected predictor X , and use it to make forecasts of a target Y**

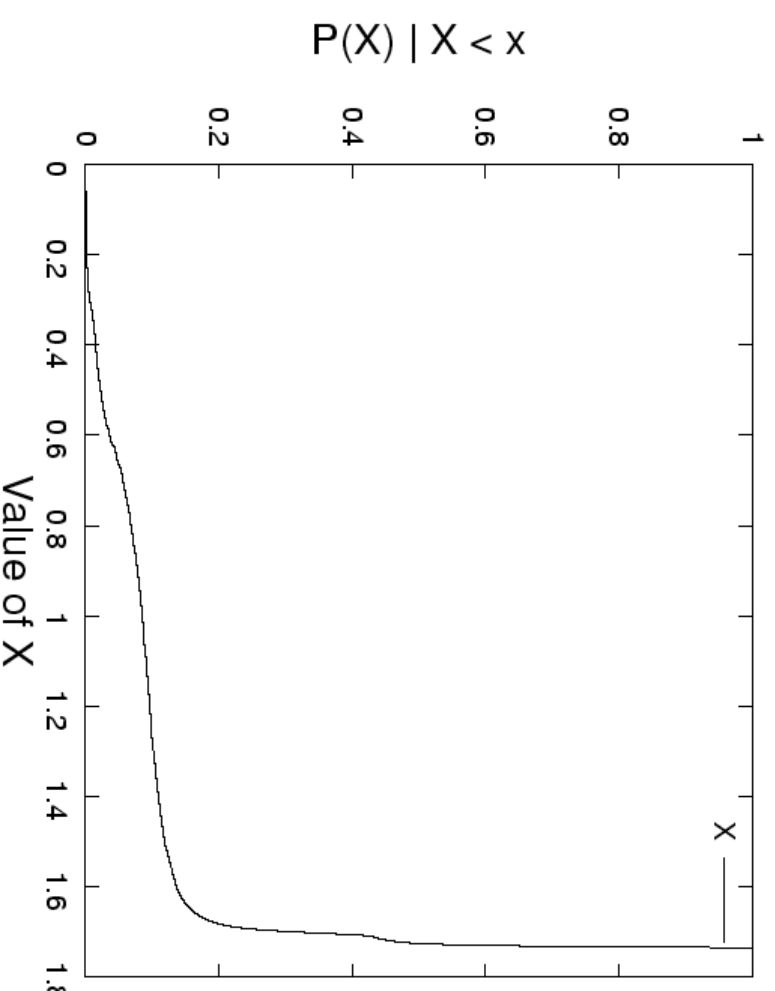
Correlation Mechanism

- Most correlation mechanisms assume normal distributions (as they deal with standard deviation, etc.)
- Network traffic does not enjoy a normal distribution
 - (see, well, lots of stuff)
- Distribution-free correlation mechanisms such as Spearman Rank correlation assume datasets of the same size

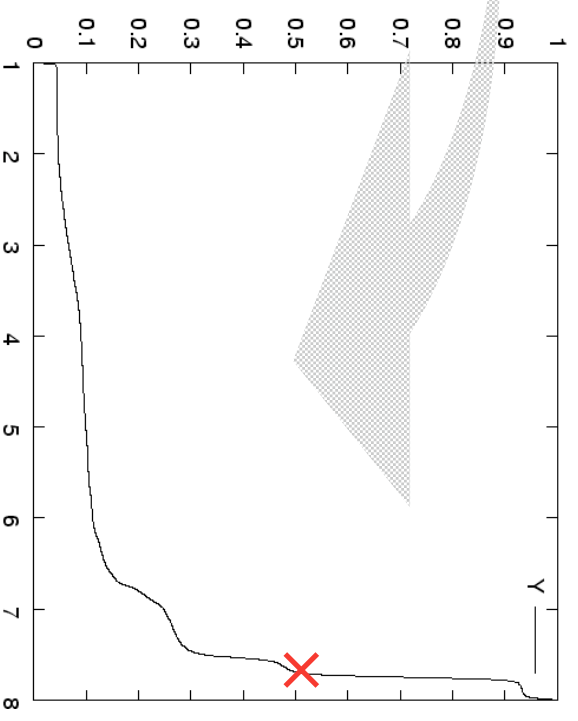
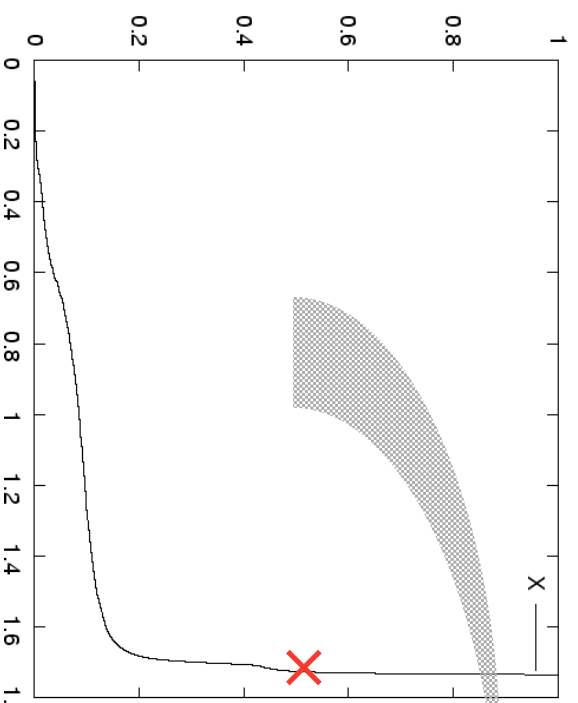
The Cumulative Distribution Function

- The empirical CDF is defined as

$$CDF_X(x) = \sum_{i=1}^{n} \frac{1}{n} \mathbb{1}_{\{x_i \leq x\}}$$



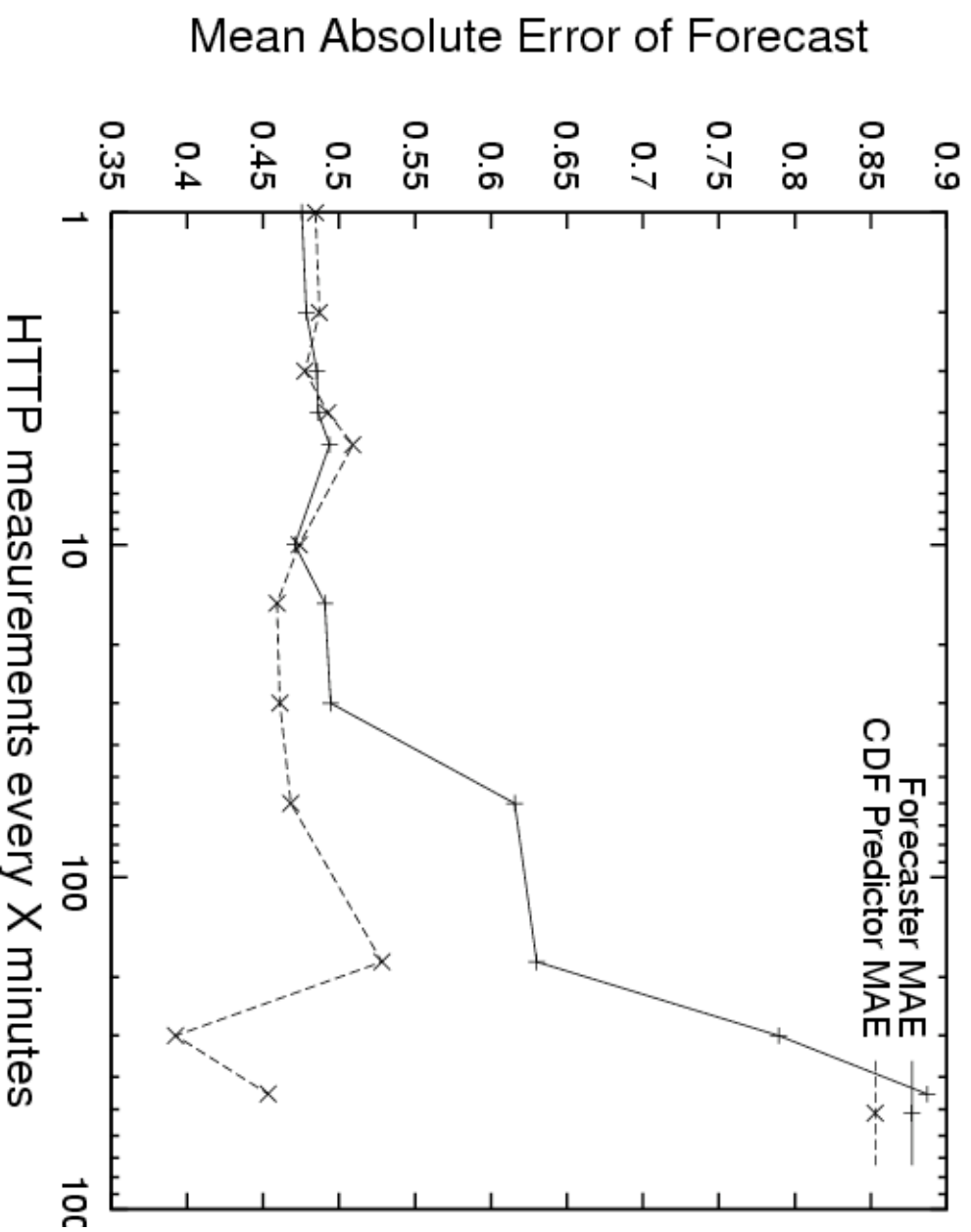
**So, our CDF correlator uses the CDF to
translate X into Y**



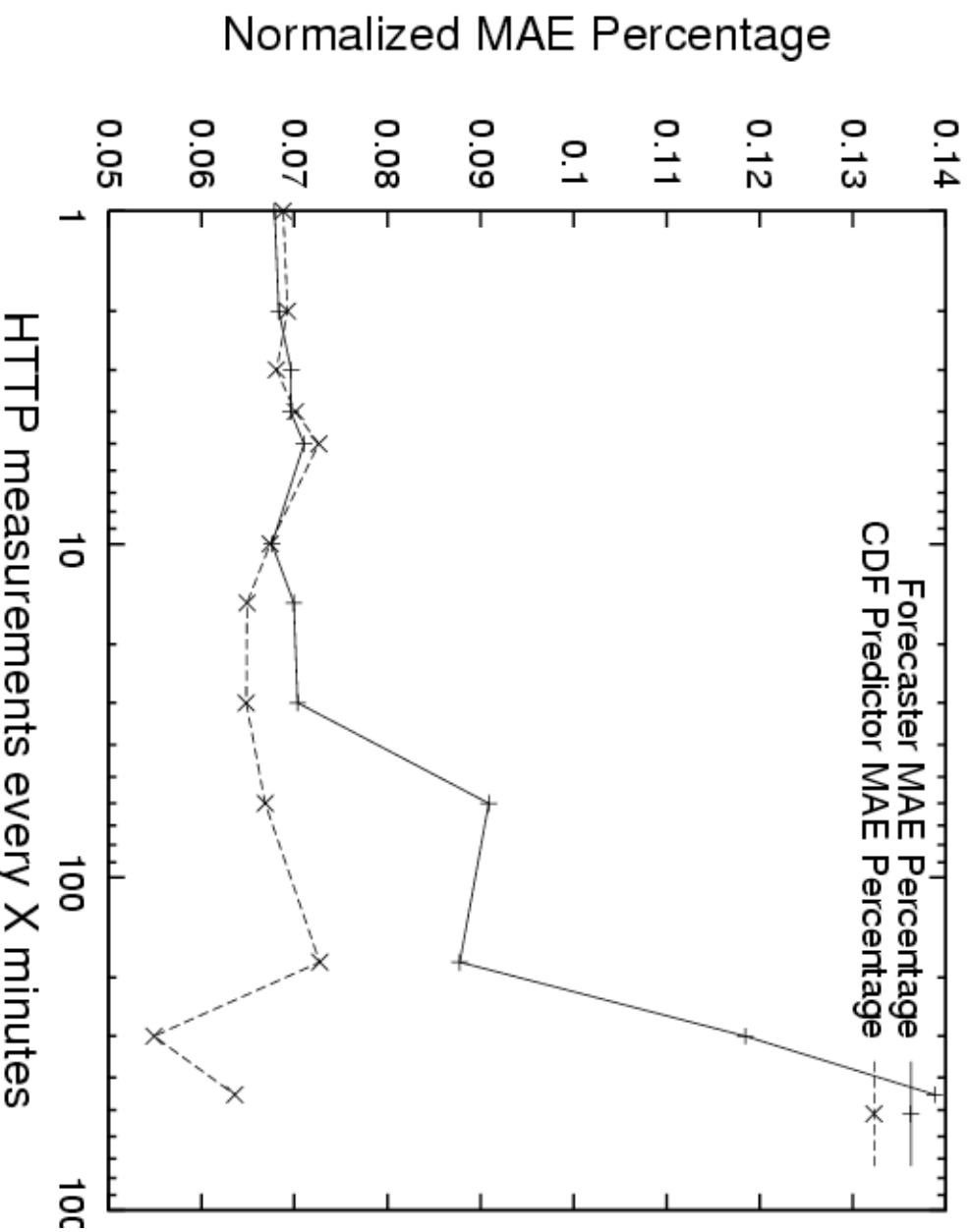
Some Terminology

- MAE – Mean Average Error
- MSE – Mean Square Error
- MNEP – Moving Normalized Error Percentage

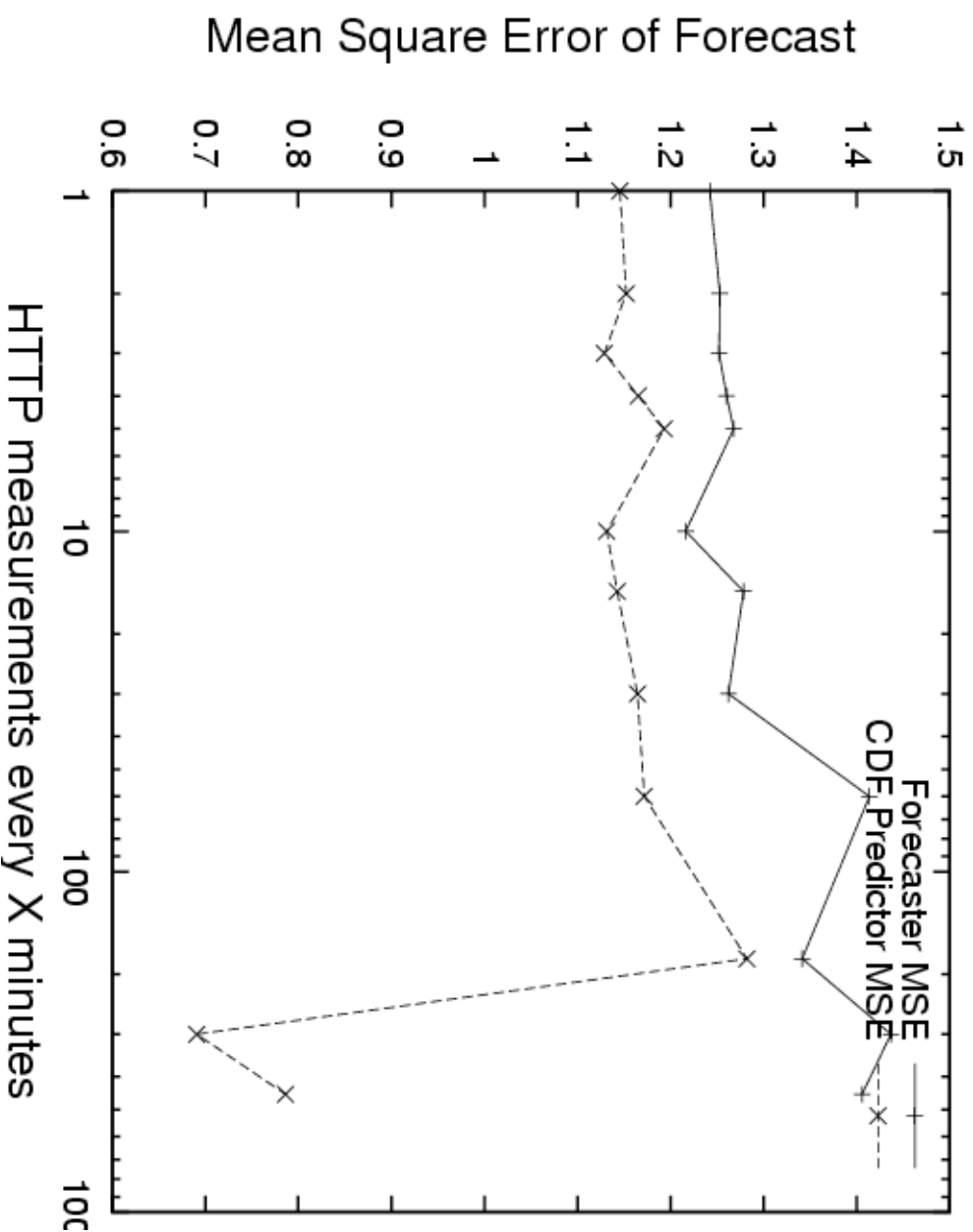
Comparison of Mean Absolute Error (MAE) between univariate and multivariate forecasts for different frequencies of HTTP measurements



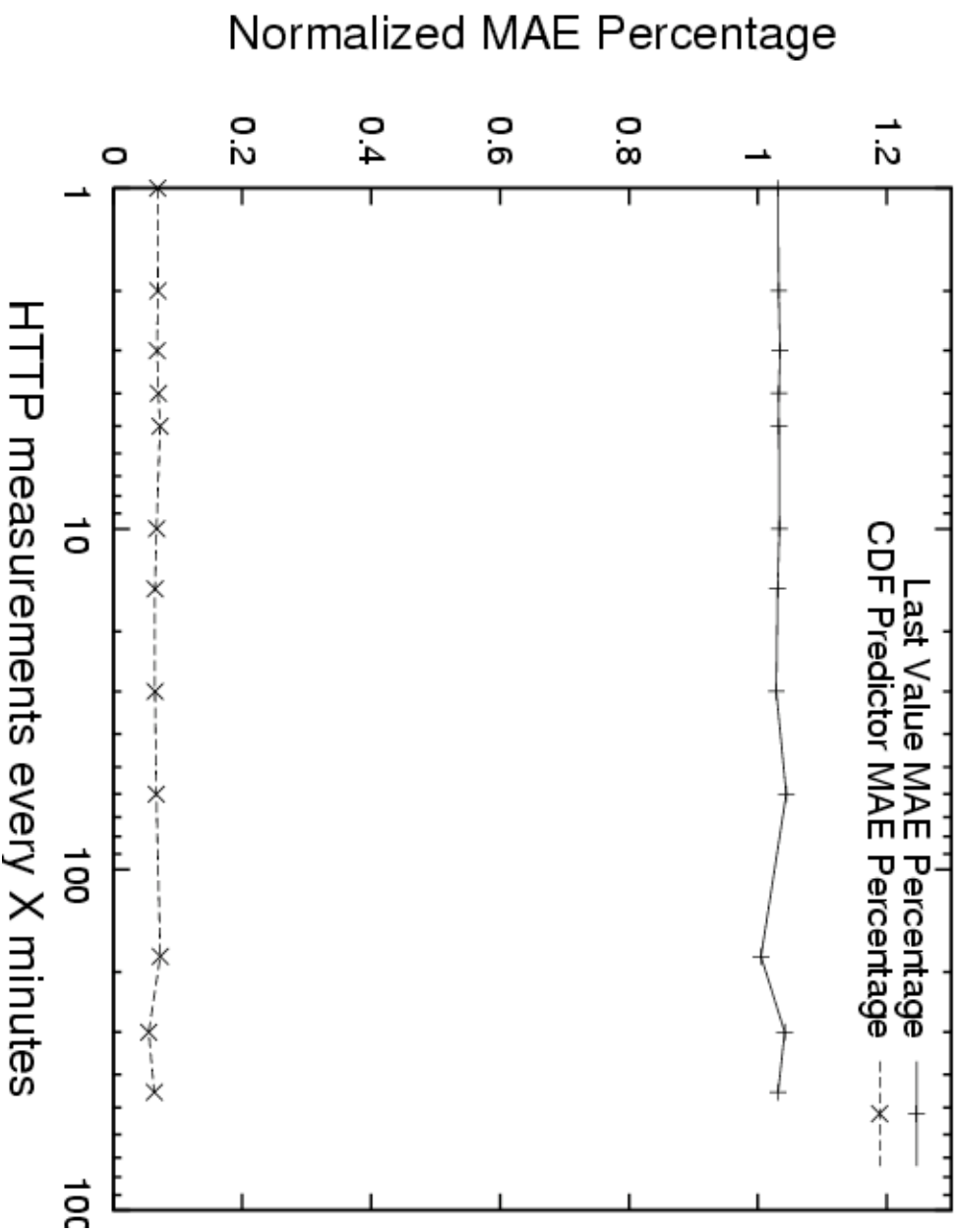
Comparison of Moving Normalized Error Percent (MNEP) of the Mean Absolute Error (MAE) between univariate and multivariate forecasts for different frequencies of HTTP



Comparison of the square root of the Mean Square Error (MSE) between univariate and multivariate forecasts for different frequencies of HTTP measurements.



Comparison of Moving Normalized Error Percent (MNEP) of the Mean Absolute Error (MAE) between "Last Value" and multivariate forecasts.



Conclusions

- We have developed a novel multivariate prediction technique
- Much yet to be done, although the CDF is proving useful in situations where a distribution is required