

# Measuring the “on-liness” of data streams

**Manfred K. Warmuth**   **Jiazhong Nie**

University of California - Santa Cruz

Dec. 10, 2015 — Nips workshop on Easy Data

Includes some earlier work with  
**Corrie Scalisi, Robert Gramacy, Scott Brandt and Ismail Ari**

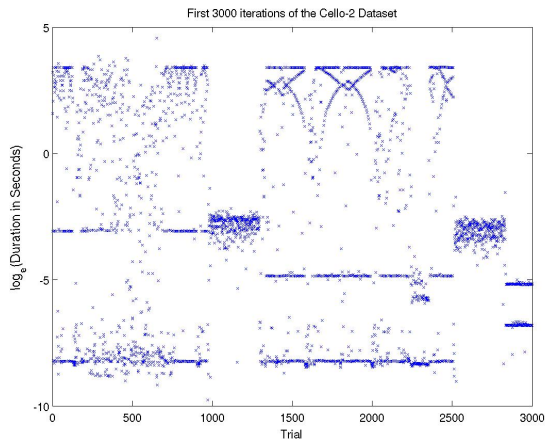
# Goals

- Design on-line algorithms in domains that are outside of the reach of theory
  
- Design good comparators that exploit the on-liness of the data

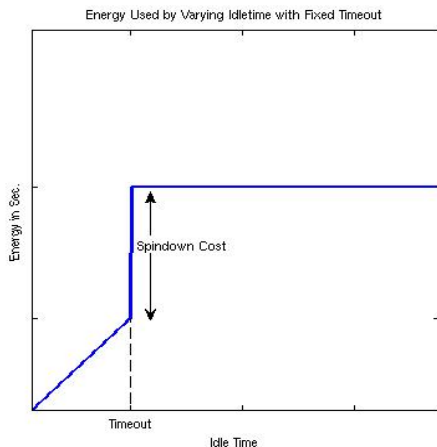
# 1. Disk spindown problem

[HLSS]

- When to spin down the disk on your laptop?
- Best time-out time/user/usage dependent



# Non-convex loss



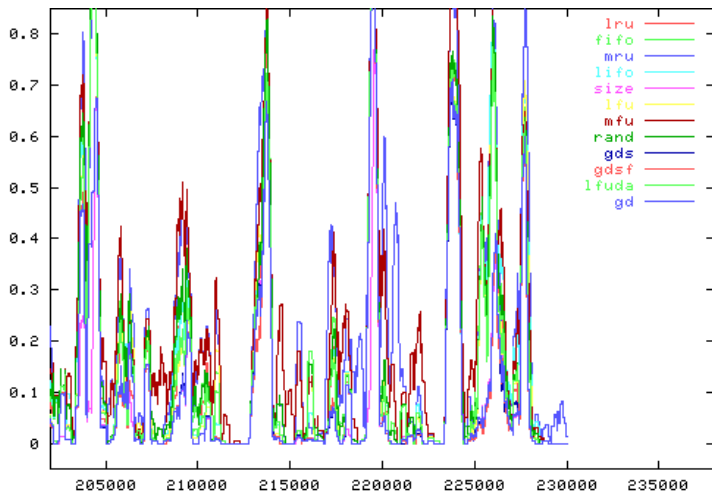
If idles times expected to be

- short, then long timeout better
- long, then short timeout better

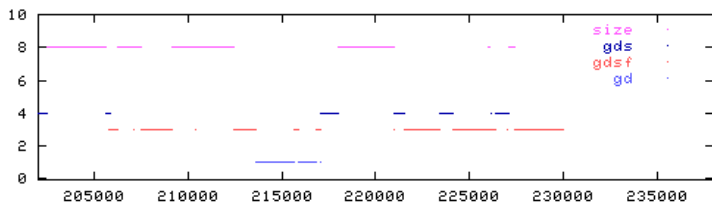
- Want to build combined caching policy from 12 base policies (our experts):

LRU, RAND, FIFO, LIFO, LFU, MFU, SIZE, GDS, GD\*, GD<sup>+</sup>, GD<sup>SF</sup>, LFUDA

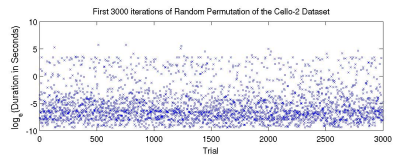
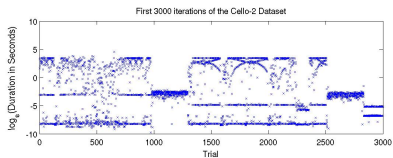
# Characteristics Vary with Time



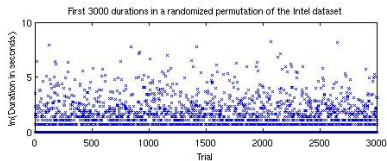
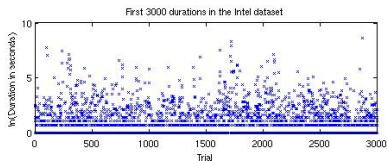
# Best Policy Varies with time



# Permuting trick for disk spindown data



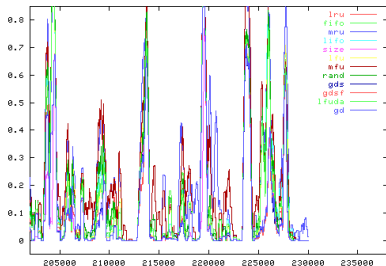
on-line :-)



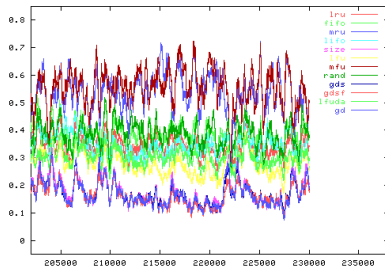
not on-line :-)



# Permuting caching data



highly on-line data



some caching policies already on-line

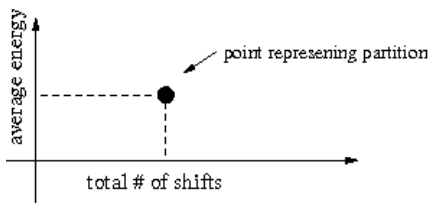
# Using a **comparators** to measure on-liness of data

## Properties

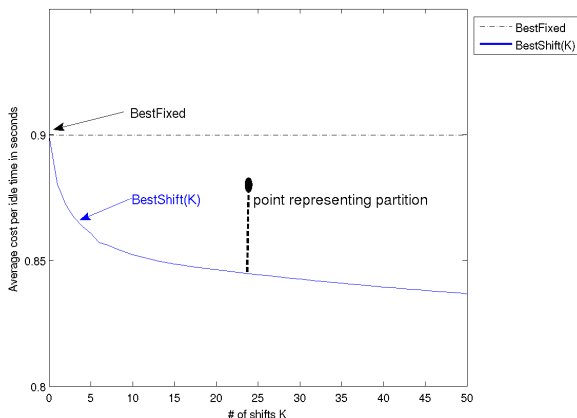
- Should exploit on-liness of data
- Might be too expensive to compute in practice, but can serve as a goal to compare against
- Might rely on information not available to the on-line algorithm

# Idea 1: Use dynamic programming to compute BestShift( $K$ ) curve

- Partition of the timeline into  $K$  segments
- BestFixed in each segment



# BestFixed( $K$ )

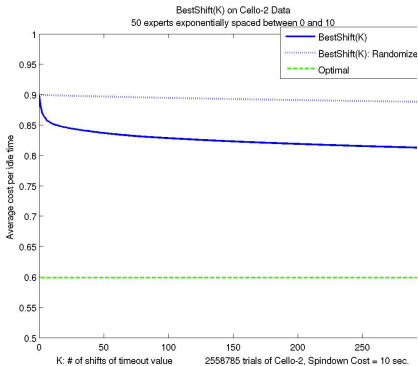


Dynamic programming:  $O(KN^2T)$

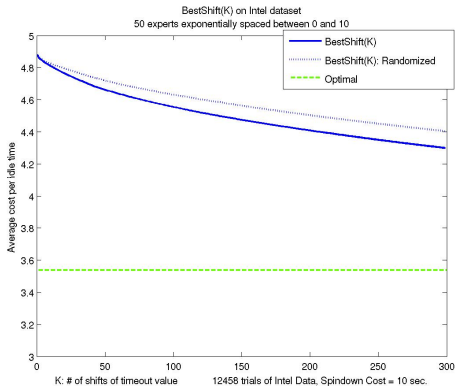
where  $K$  # of partitions,  $N$  # of discrete idle times,  $T$  # of trials

[H]

# BestShift curves



on-line



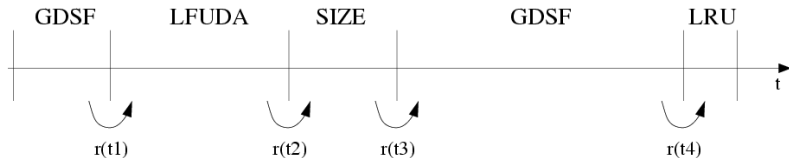
not on-line

# Comparators for caching

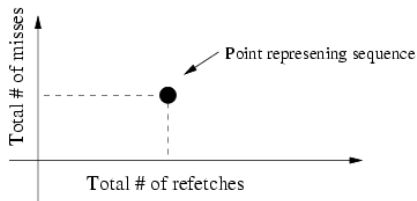
- **BestFixed**: a **posteriori best** of 12 policies on entire request stream
- **BestRefetching**( $R$ ):  
minimum number of misses with at most  $R$  refetches  
in any sequence of switching policies

# Refetches & Policy Switches

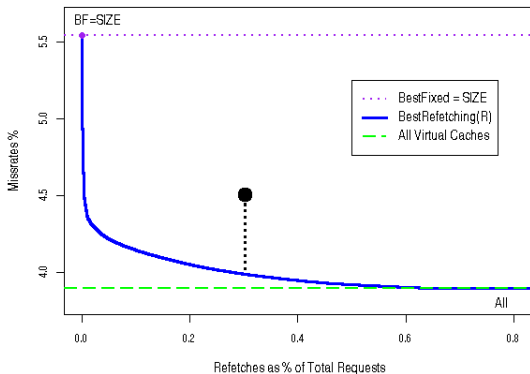
**Comparator:** All sequences of the form



We plot miss rate v.s. refetches:



# BestRefetching( $R$ )



Dynamic programming:  $O(RN^2T)$

[H]

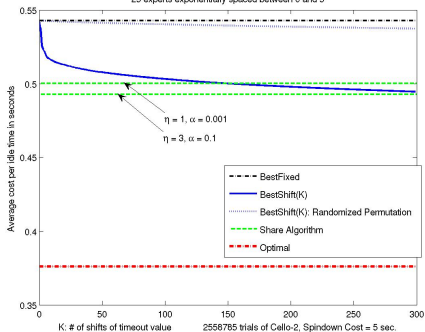


# Our theoretically sound algorithms become heuristics

- Use loss and share updates on non-convex losses
- Build a merged cache that does not correspond to the mixture

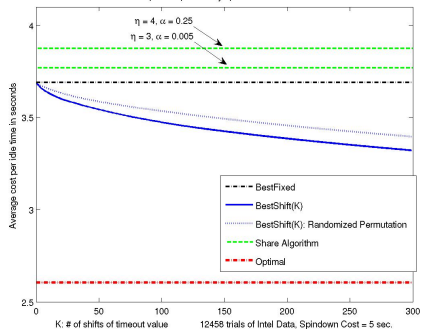
# Spindown results

BestShift(K) and Share Algorithm on Cello-2 dataset  
25 experts exponentially spaced between 0 and 5



on-line :-)

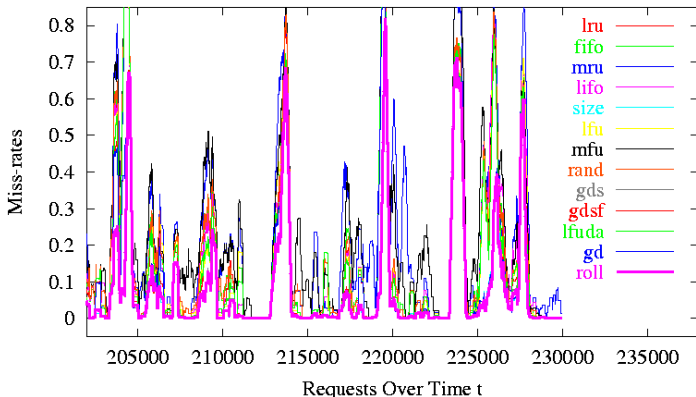
BestShift(K) and Share Algorithm on Intel dataset  
25 experts exponentially spaced between 0 and 5



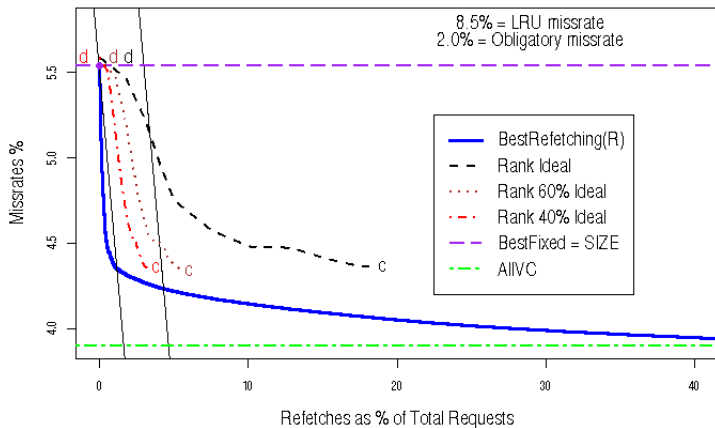
not on-line :-)

# Caching - we “Tracks” best policy

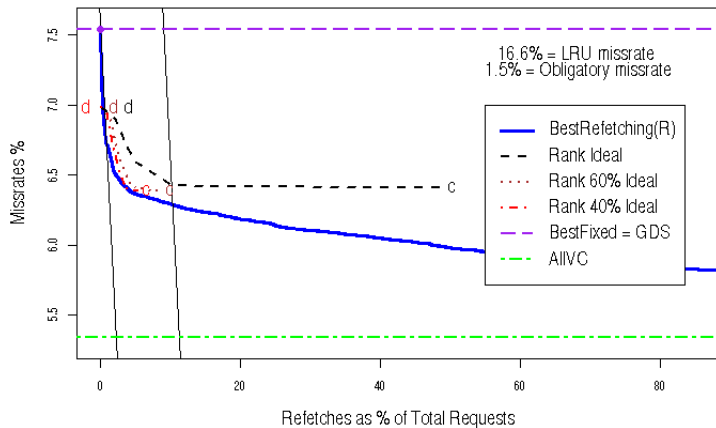
Miss-rates under FSUP with Master



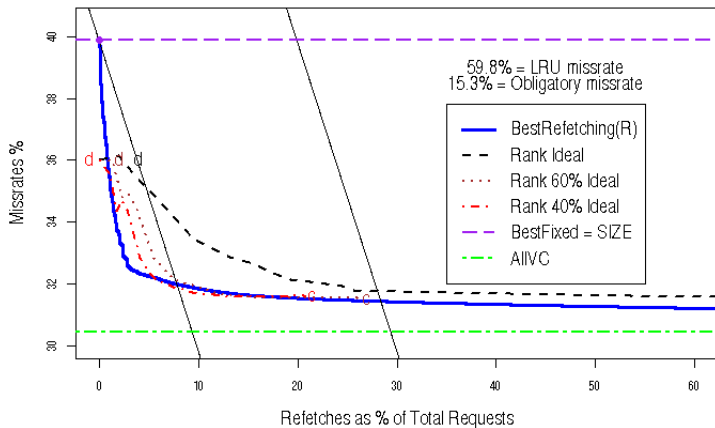
## WWk Master and Comparator Missrates



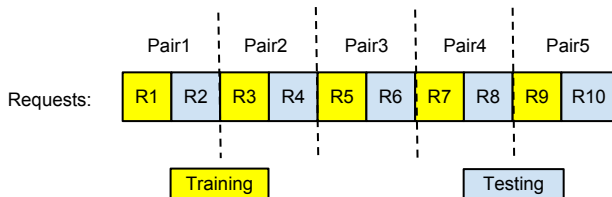
## UMo Master and Comparator Missrates



SMoLRU Master and Comparator Missrates



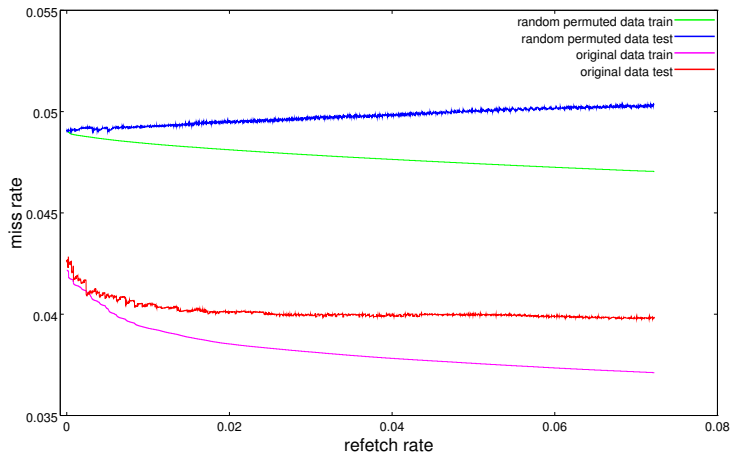
## Idea 2: Split into even/odd requests



- Best partition based on training set
- Performance based on test set

# Miss Rate of Testing Requests

No overfitting to random data: testing miss rate goes up immediately





# Upshot!

Don't be afraid  
to use your algorithms  
as heuristics in domains  
where the theory breaks down