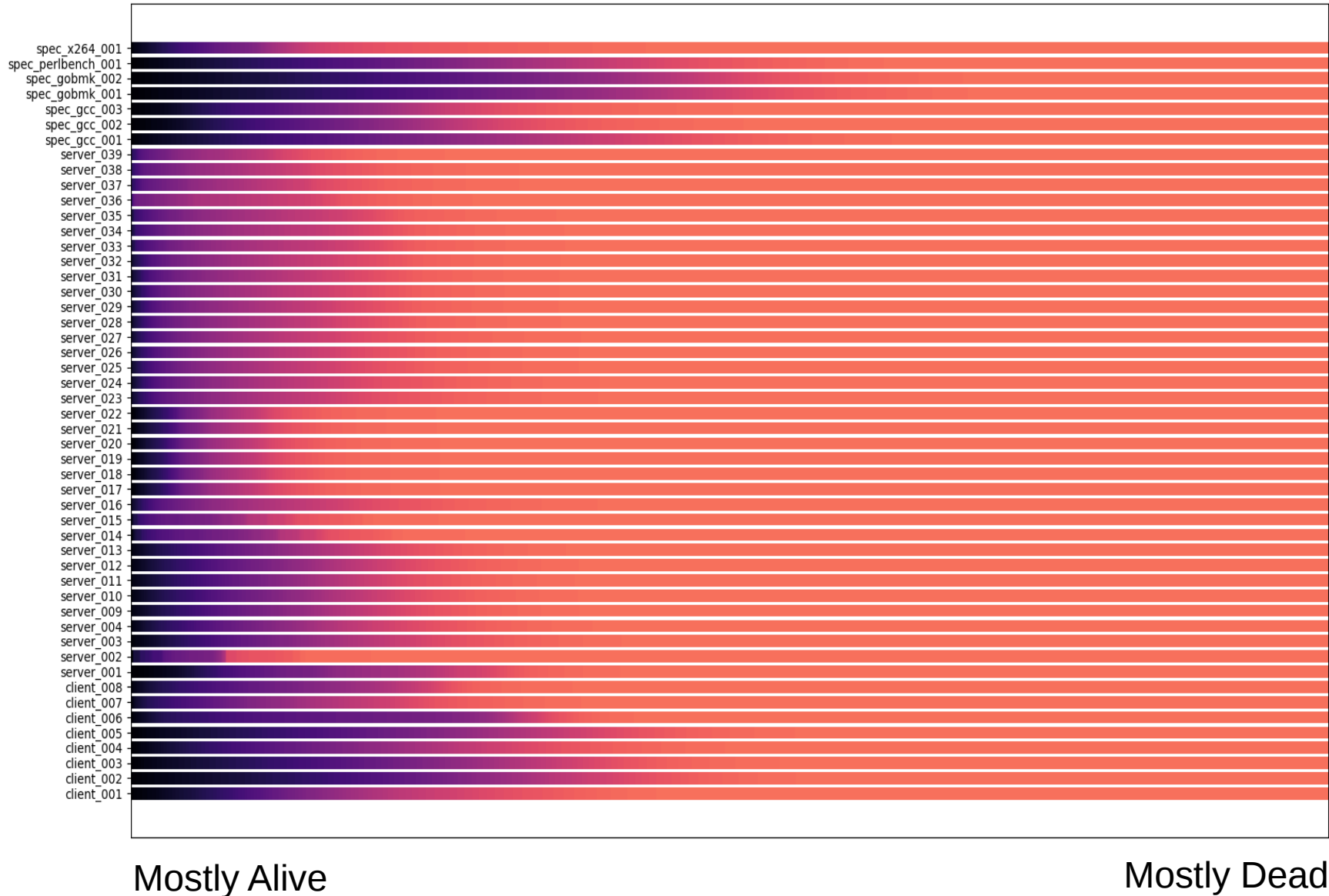


Temporal Ancestry Prefetching

Nathan Gober, Gino Chacon,
Daniel Jiménez, Paul V. Gratz
Texas A&M University



Block Dead Times in L1I

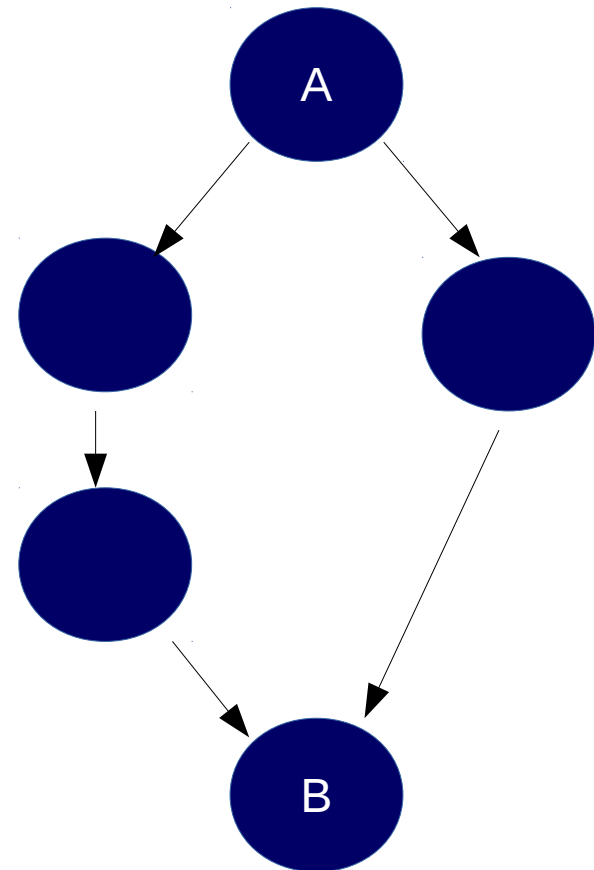


Outline

- Background
 - Ancestry Prediction
 - Temporal Prefetching
- Design
- Results
- Conclusion

Ancestry prediction

- Many paths of varying lengths from A to B
- Blocks in I-cache are long-lived
- Short histories miss the connection between A and B



Temporal Prefetching

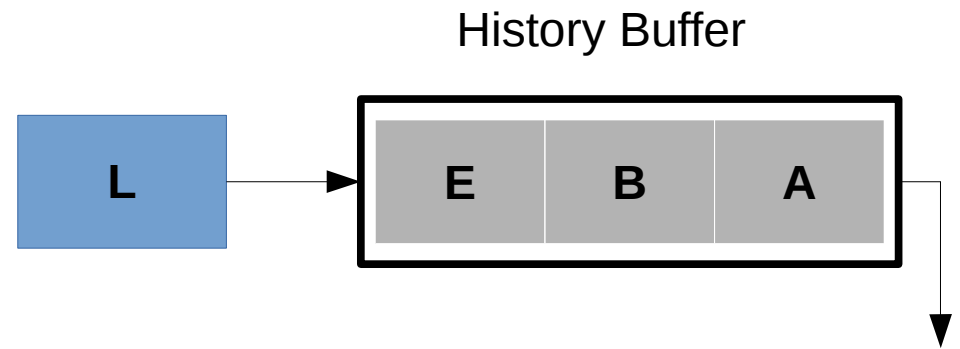
- Replay old misses on cache access
- Large metadata requirements
 - State of the art uses off-chip memory [Wenisch HPCA'09, Bakhshalipour HPCA'18]
 - Structural address space? [Wu ISCA'19]
- IPC1 gives large hardware budget

Outline

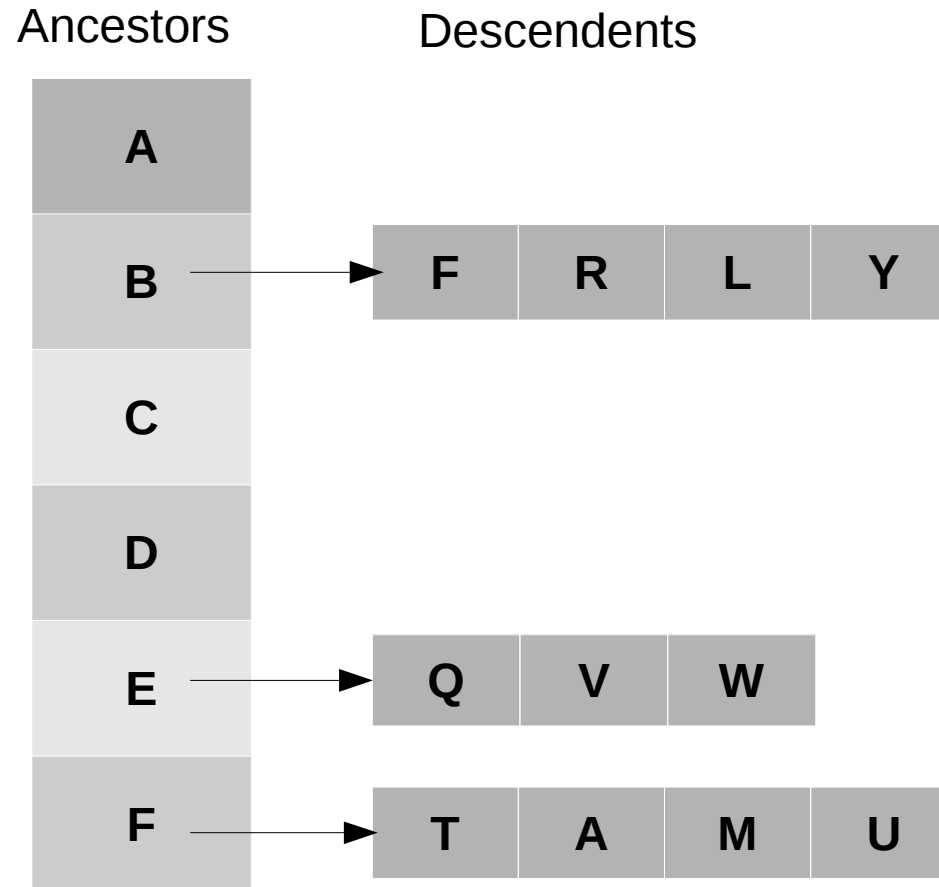
- Background
 - Ancestry Prediction
 - Temporal Prefetching
- Design
- Results
- Conclusion

TAP Model (Training)

- Maintain path history
- De-duplicated
- A single PC has many histories

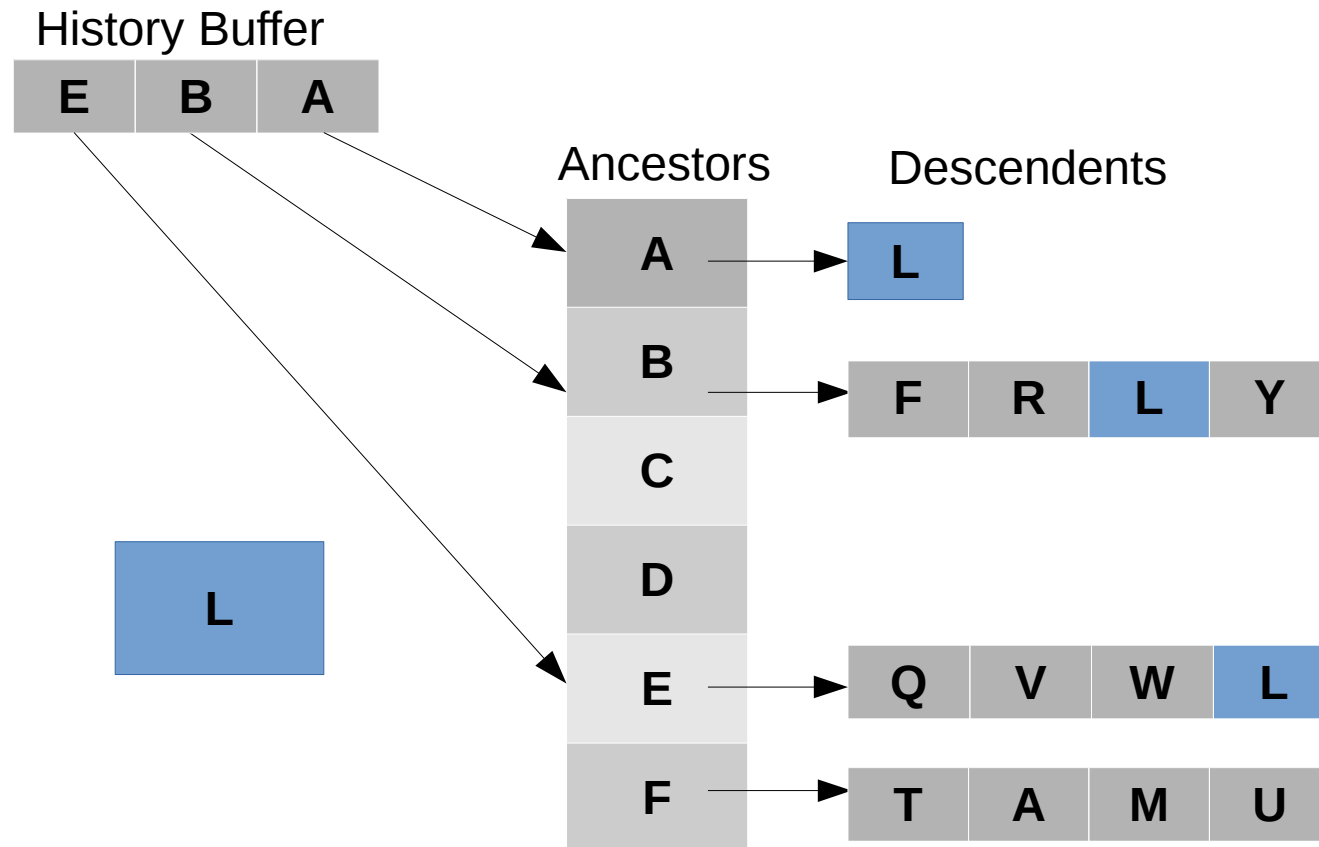


TAP Model (Training)



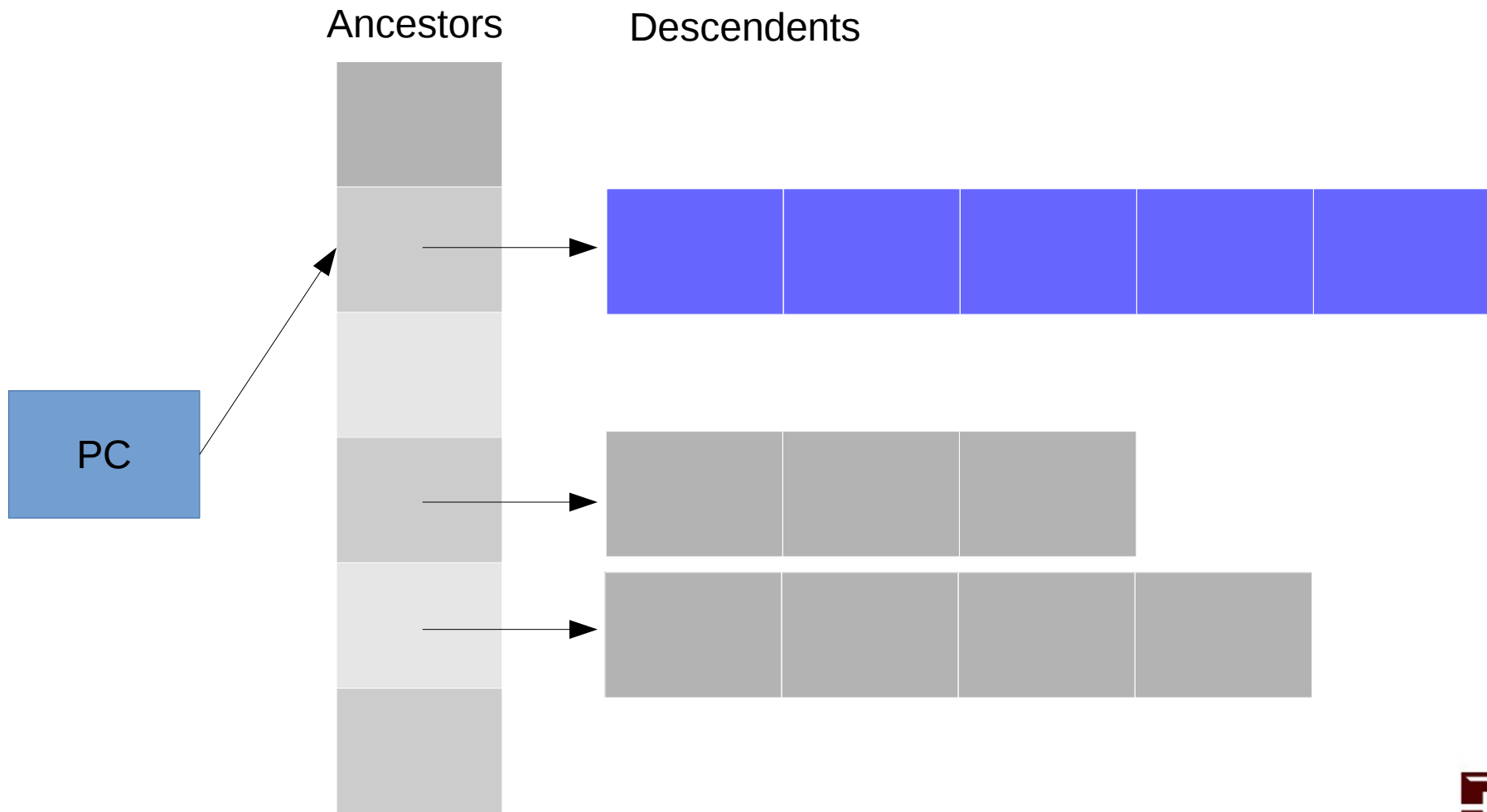
- Descendants are captured misses
- The path is not recorded, only the eventual predicted miss

TAP Model (Training)



- Increment weights on cache access
- Decrement weights on eviction if prefetch not useful
 - Invalidate instead if weight was zero.

TAP Model (Prediction)

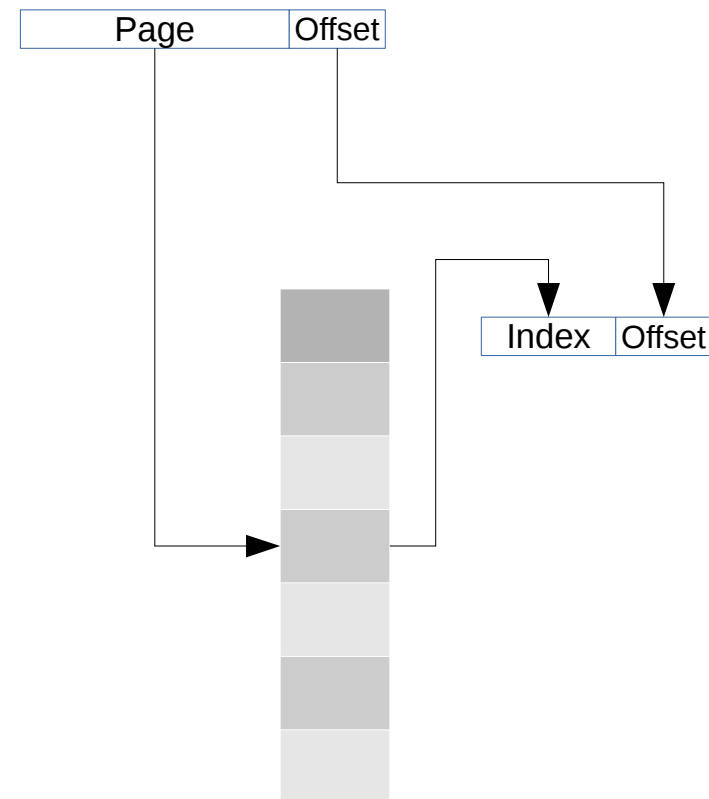


Prefetch Filtering

- Long descendency lists produce large numbers of prefetch hits.
- Filter prefetches by tag check.
- With shadow cache of 12-bit partial tags, filtering is 99% accurate.

Page Compression

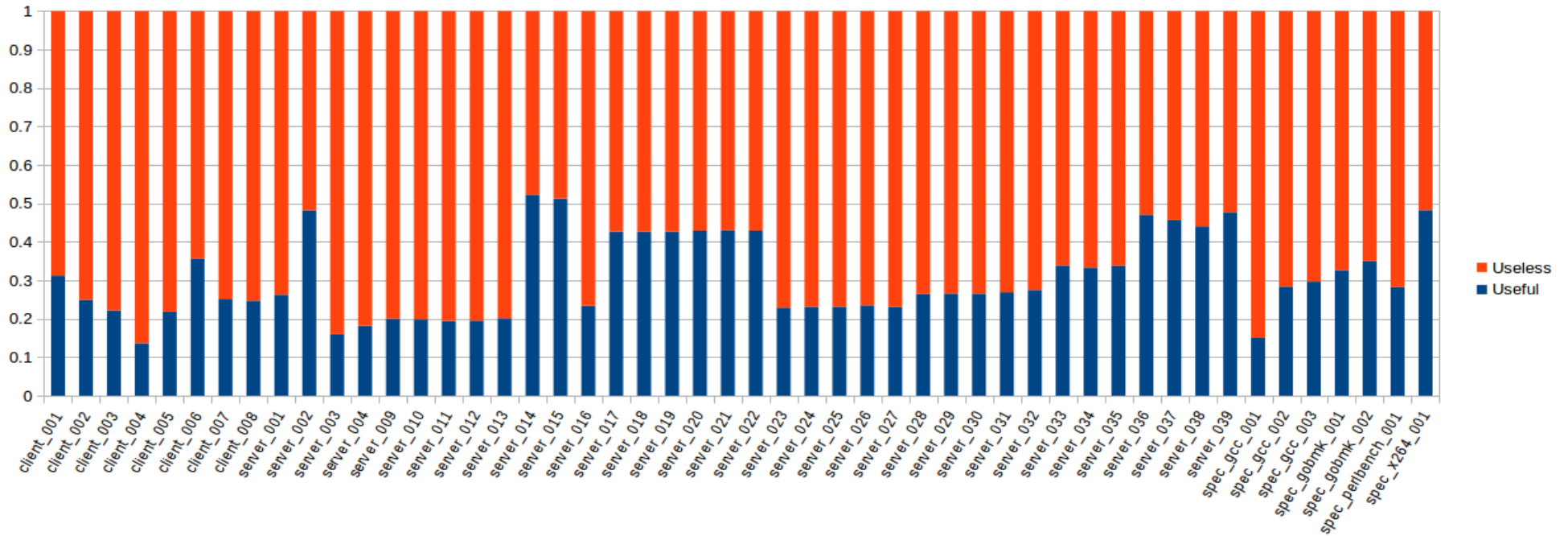
- Much metadata space is wasted on storing a few pages.
- Use buffer index as a proxy for page number.
- NRU replacement



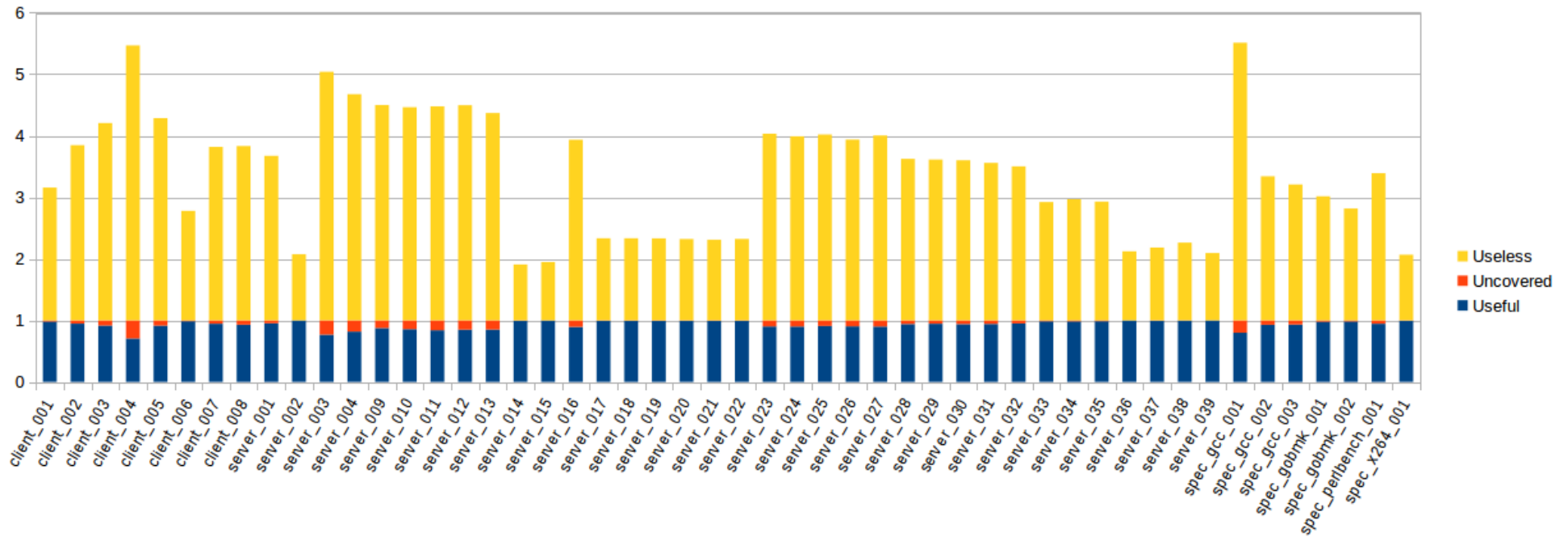
Outline

- Background
 - Ancestry Prediction
 - Temporal Prefetching
- Design
- Results
- Conclusion

Accuracy



Coverage



Conclusion

- Tracking cache evictions is useful
- Prefetch sequencing is unimportant
- Future Work
 - Replacement in L1I
 - Stricter throttling on poor-performing workloads