# Flash as a Cache: Challenges & Opportunities
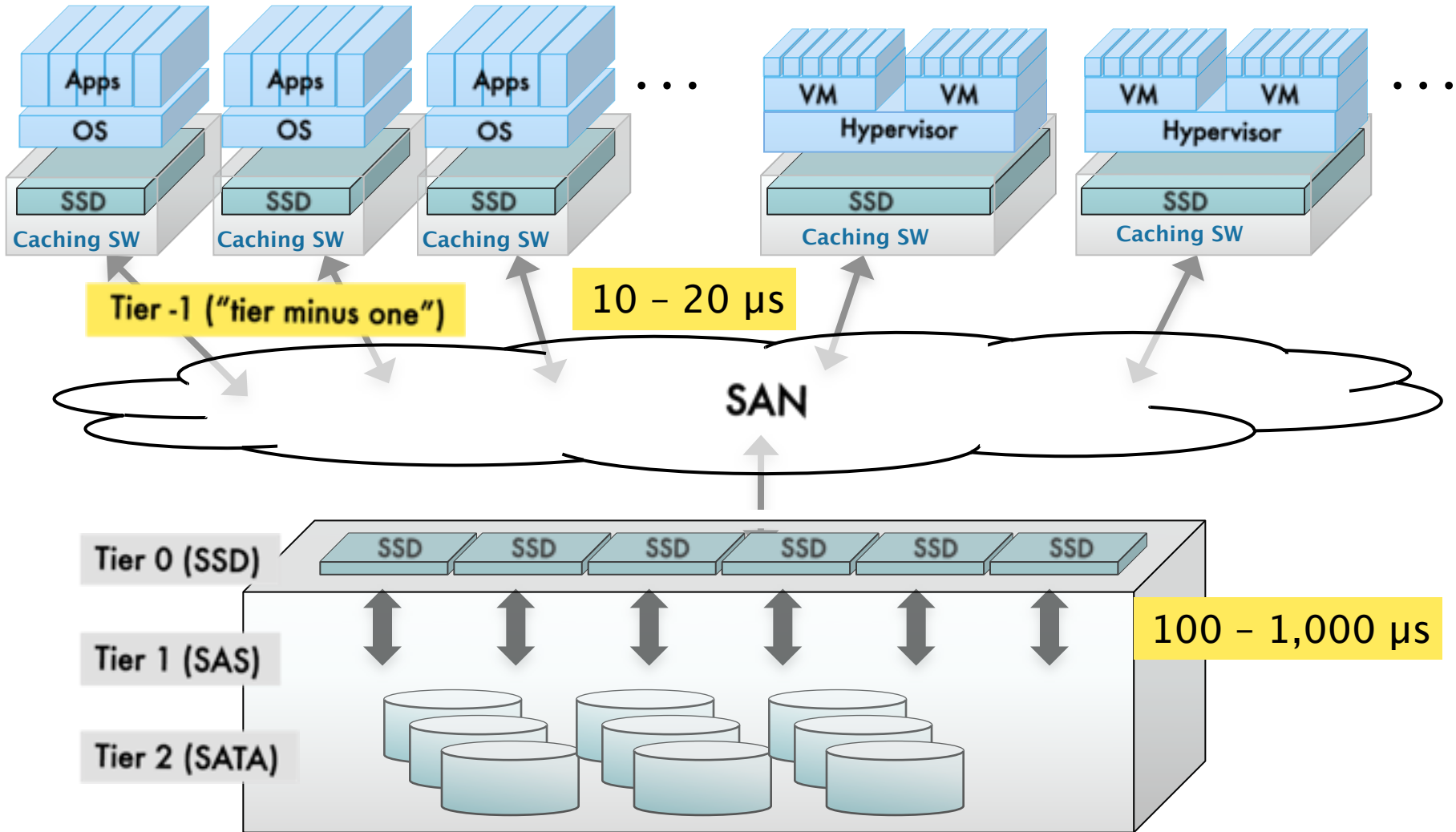
## Ted Sanford
## CEO, FlashSoft

# Agenda

- **Technical Choices & Challenges**
  - What is "tier minus one" caching?
  - How can it be implemented?
  - When is it useful?
  - Key design considerations
- **Use Cases & Performance Impact**
- **Architectural Considerations**
  - Accelerating a stand–alone server
  - Distributed cache for server cluster
  - Extending System memory
  - Multi–level solid state cache design
  - Managing IO resources

# What is T −1 caching?

# When is T –1 caching useful?

- IO must be a significant bottleneck
  - CPU not fully utilized
  - Sufficient memory available (esp for VMs)
- Data must have hot spots
  - The 80–20 rule
  - The 90–90 rule – hot spots are always changing
  - Goal: matching the economic value of the data to the appropriate IO resource
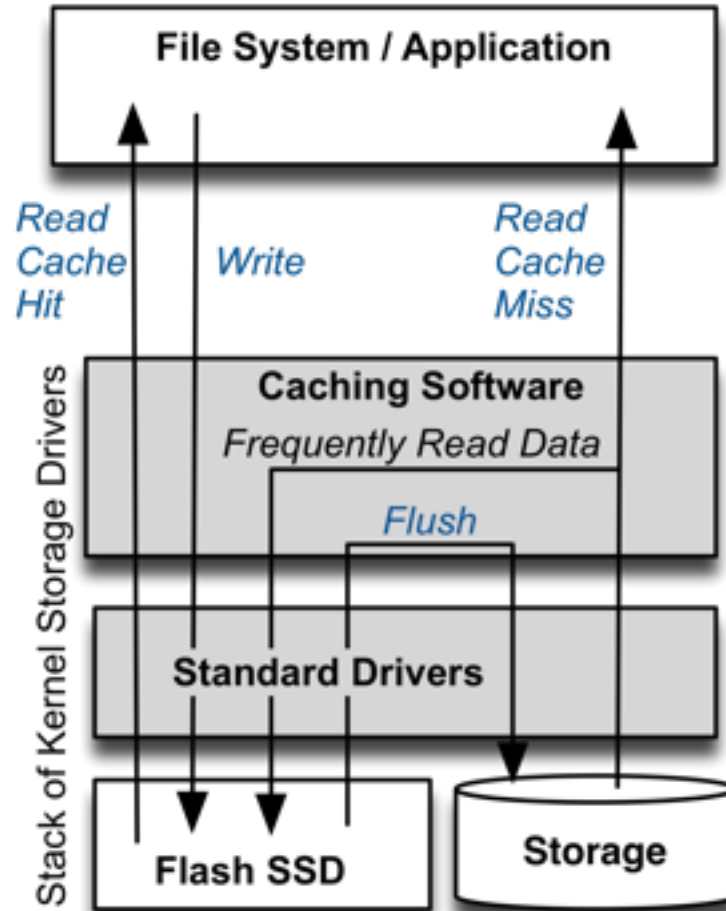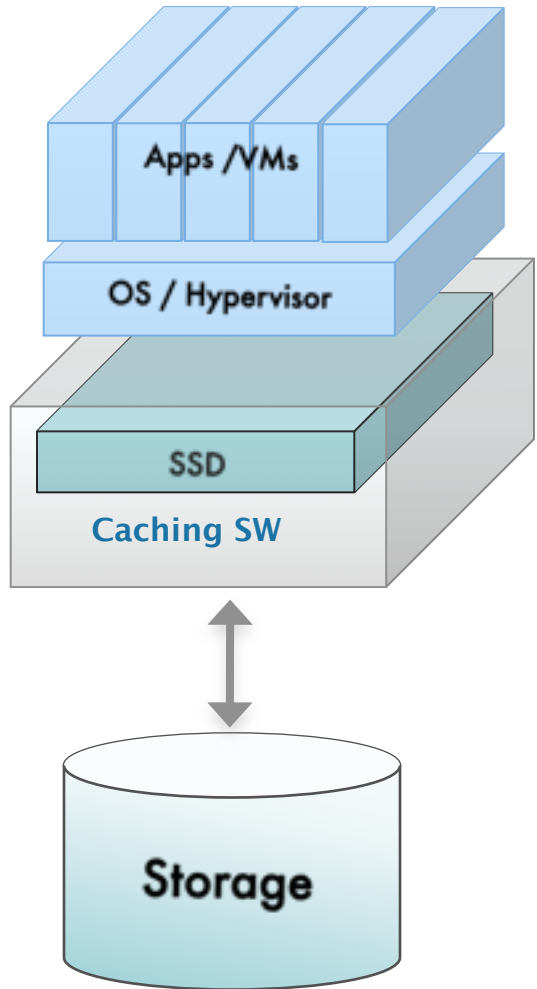
# Key Design Considerations

- Designing cache optimized for flash architecture
  - Traditional buffer cache design
  - Circular buffer, log-structured cache
- Minimizing server overhead
  - Memory
  - CPU
- Read-only vs. Read-Write
- Accommodating diversity of SSD designs
  - Size
  - Speed
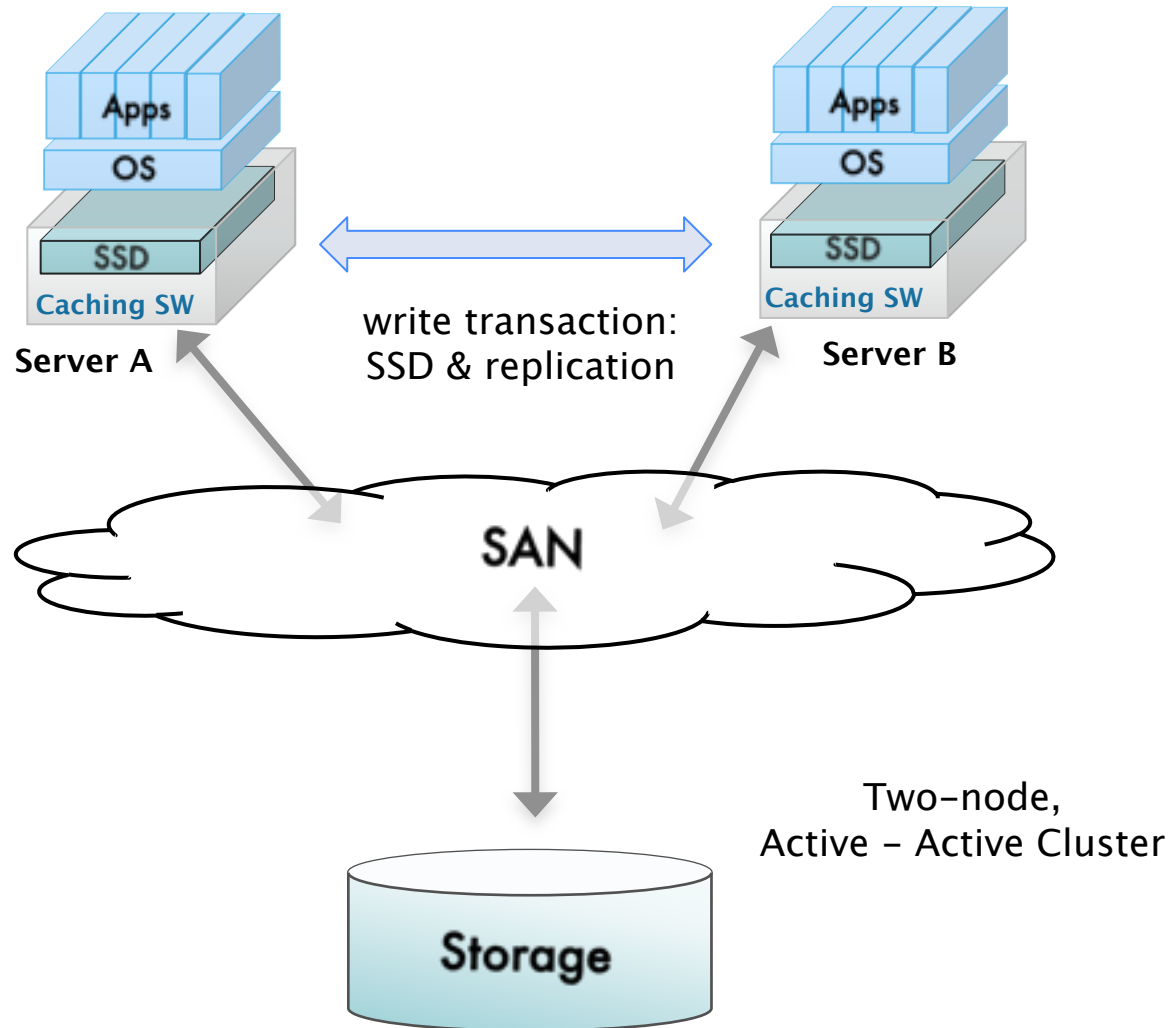  - Block architecture
  - Controller

# Use Cases

- Accelerating applications on stand-alone and cluster servers
  - Windows
  - Linux
  - Virtualization
- Applications most likely to benefit
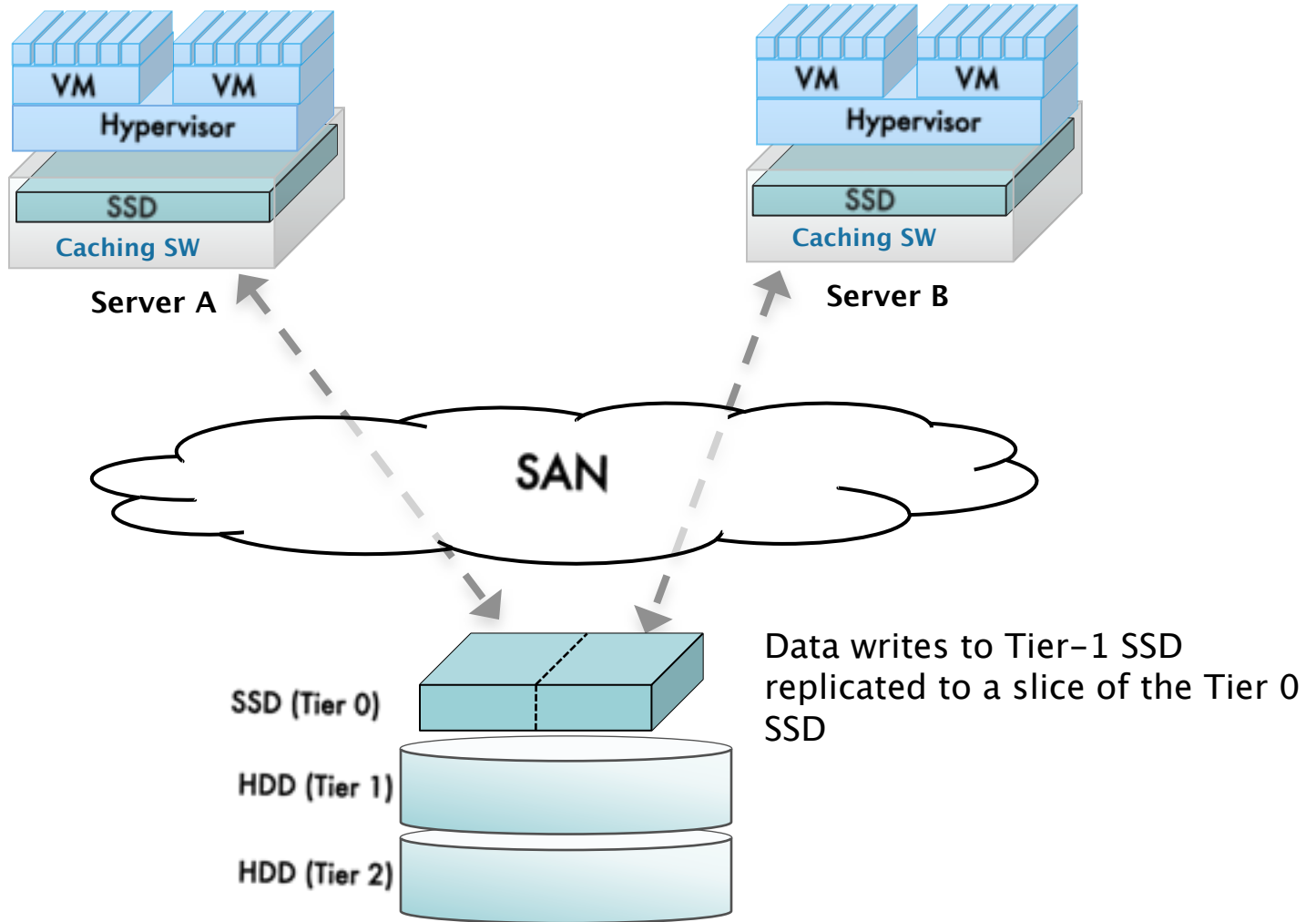  - Databases: Analytics, BI, ERP
  - Websites: Social, Search, SaaS, E2.0

# Caching in a stand-alone server
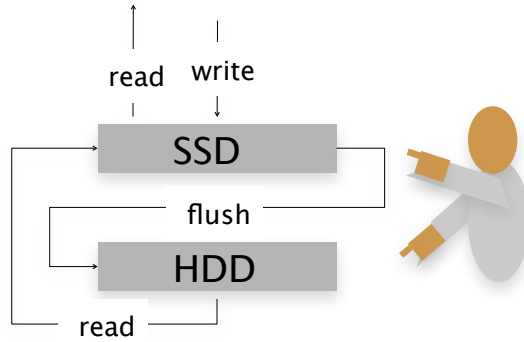
# Cluster: Horizontal Replication



**Server A**

write transaction:
SSD & replication

**Server B**

SAN

Two-node,
Active – Active Cluster

Storage

# Cluster: Vertical Replication API



Data writes to Tier-1 SSD replicated to a slice of the Tier 0 SSD

# Multi-level solid state cache design

**Today:**

read    write

SSD

flush

HDD

read

**Two Tier caching**
- allocates data to either SSD or HDD

**Future:**

read    write
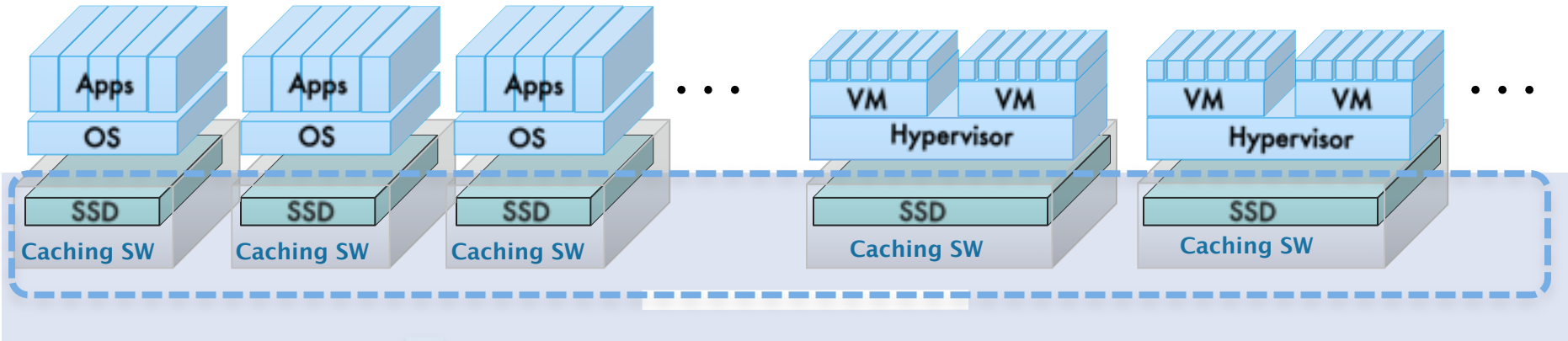
DRAM

PCIe SSD

SAS/SATA SSD

HDD

**Multi-tier caching**
- allocates data to DRAM, PCIe SSD, SAS/SATA SSD, or HDD
- Extensible to NVMe, PCM

# Managing IO resources

# Thank You