# Elasticity

# No storage elasticity

# Not all data requires durability



**Significant ephemeral data**

What if files could **contract** and **expand**?

# The motif abstraction

## A **motif** is **code** to generate the **data** in a file

- **expand()** obtains the raw bytes of the file
  - Run **computations** (e.g. **compression**)
  - **Fetch data** across a network
  - Operate over **other files** on the FS

- **contract()** deletes raw bytes, retains motif code
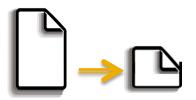
| Motifs can be recursive | Motifs can have circular dependencies | Files can have multiple motifs |

# Example *network storage* motif

```python
class SCPMotif(object):
  def expand(self, fname, meta=None):
    p = fname.bypass()
    os.popen('scp \
        fileserver1:storage%s "%s"' % (p,p))

  def contract(self, fname, meta=None):
    p = fname.bypass()
    os.popen('ssh fileserver1\
      "mkdir -p storage%s"' %\
      os.path.dirname(p))
    if os.popen('scp "%s"\
        fileserver1:storage%s' % (p,p)) == 0:
      open(p, 'w').close()
```

The harmonium file system

VM

APPS

harmonium FS

20GB

Virtual disk          250GB

Limit on total size

Large address space

Hypervisor

# The harmonium file system

VM

APPS

harmonium FS **10GB**

20GB

**Virtual disk** 250GB

Conductor

**Hypervisor**

The harmonium file system

VM

Conductor

Hypervisor

**Which files** should be contracted or expanded?

**What interface** can the conductor use?

# What files to contract/expand?

harmonium FS     $x$ **GB limit**

**When will the file be accessed next?**

**How much time will it take to expand?**

- **Feed info into optimization mechanism**
  - **Choose files with minimum total expansion latency s.t. contraction saves sufficient space.**
  - **Most competitive algorithm: Greedily choose files on LRU list to maximize ratio of space savings to expansion latency**

# Conductor interface

**VM**

**harmonium FS**

20GB

**Virtual d...**

10GB

**TRIM**

**Conductor**

**Hypervisor**

# Evaluation

- **User-space FUSE prototype**
  - **Conductor: Python, via UNIX sockets**

- **Workload**
  - **Set of 54,000 patch files applied in chronological order to the Linux kernel**

- **Motif**
  - **Network storage via scp**

- **Measure latency to access first byte**

Top chart — Time-to-first-byte latency (ms) vs Trace time (%), with regions labeled: Start, Contraction, Managed, Expansion, End.

Bottom chart — File system size (MB) vs Trace time (%), with regions labeled: Start, Contraction, Managed, Expansion, End.

Annotations: Use ≤450MB, Use ≤650MB, No size constraints, Size constraint lifted

# Conclusion

Elastic performance/capacity trade-off for storage in VMs

## Problem

Applications store **ephemeral** data on secondary storage

But storage stacks provide **durability** for every file

## harmonium

Associate *motifs* with every file, allowing reconstruction

**Contract/expand** files to minimize access latency

## Results

Best algorithm: LRU **greedily** maximizing space/latency

Fully functional **FUSE** prototype

# Extra slides

# Theoretical formulation

- **0-1 Knapsack (NP-complete)**
  - *$S$* = space needed to save
  - $e_i$ = expected expansion latency for file *$i$*
  - $s_i$ = expected storage savings for file *$I$*

$$\min \quad \sum_{i=1}^{n} e_i x_i$$

$$\text{s.t.} \quad \sum_{i=1}^{n} s_i x_i \geq S,$$

$$x_i \in \{0, 1\} \quad , \quad 1 \leq i \leq n$$

- **Reduction:** Choose the files *not* to include
- **APX-knapsack** takes *O(nW)* time, prohibitive

# Related work

- **Trade-off between storage footprint and performance**
  - **Usually in distributed settings**
  - **Sierra (EuroSys 2011), Rabbit (SOCC 2011), Springfs (FAST 2014), …**

- **These systems maintain 1 to N copies of each file**

- harmonium, however, maintains "0 to 1" copies of each file

# But isn't this just …

- **A compression file system?**
- **A glorified cache?**
- **A de-duplication system?**

  - **Harmonium can support arbitrary motifs:**
    - **Compression of rarely accessed files**
    - **Remote network storage [scp/rsync/nfs/…]**
    - **Pull from repositories [git/svn/…]**
    - **Re-wget data in Downloads folder**
    - **Resume torrent download [remove partial files]**
    - **System packages [retrieve from apt/debian/…]**
    - **Regeneration of data sets [ala Nectar]**
    - **…**

# Security concerns

- **Motifs are really arbitrary code**
  - **Can cause system to hang, crash, corrupt data or consume resources wastefully**
- **Our current implementation is vulnerable**
  - **Motifs execute within the same process as the FS**
  - **Isolation by virtualization too coarse-grained**
  - **Sandboxing great for security, may be slow**

- **Ongoing work: require authorization**
  - **Users specifically approve running of motifs generated by those of lesser privilege or fewer capabilities**

# Computation vs. storage

- *"But isn't computation more expensive than storage?"*
- **Underlying principle of our work:**

   **Computation, Network and Storage are fungible**

- **Harmonium allows use of Computation or Network when Storage is scarce**
- **Other parts of the trade-off interesting as well!**