

Bolt: Building a distributed ndarray

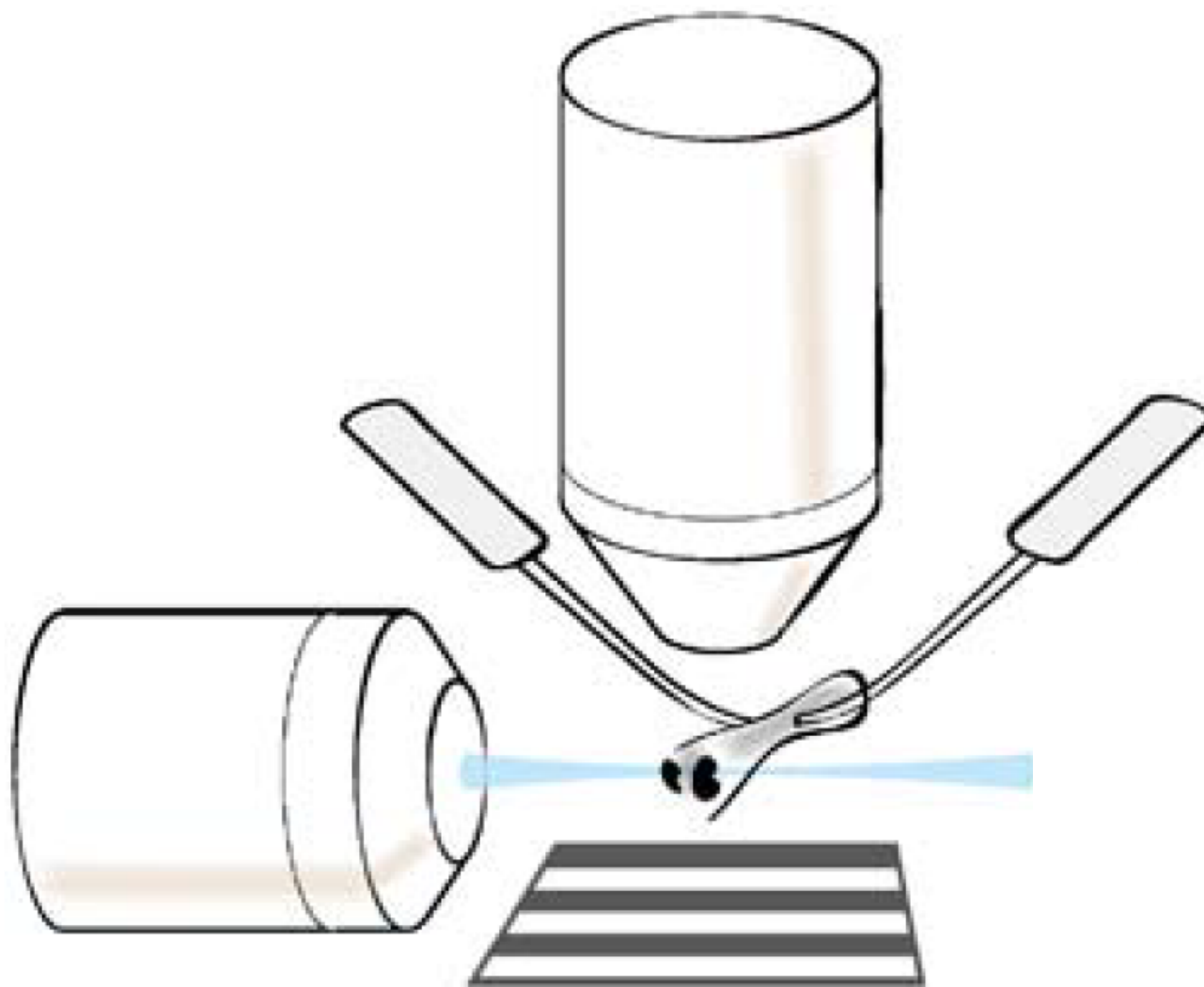
Jason Wittenbach

Janelia Research Campus (HHMI)

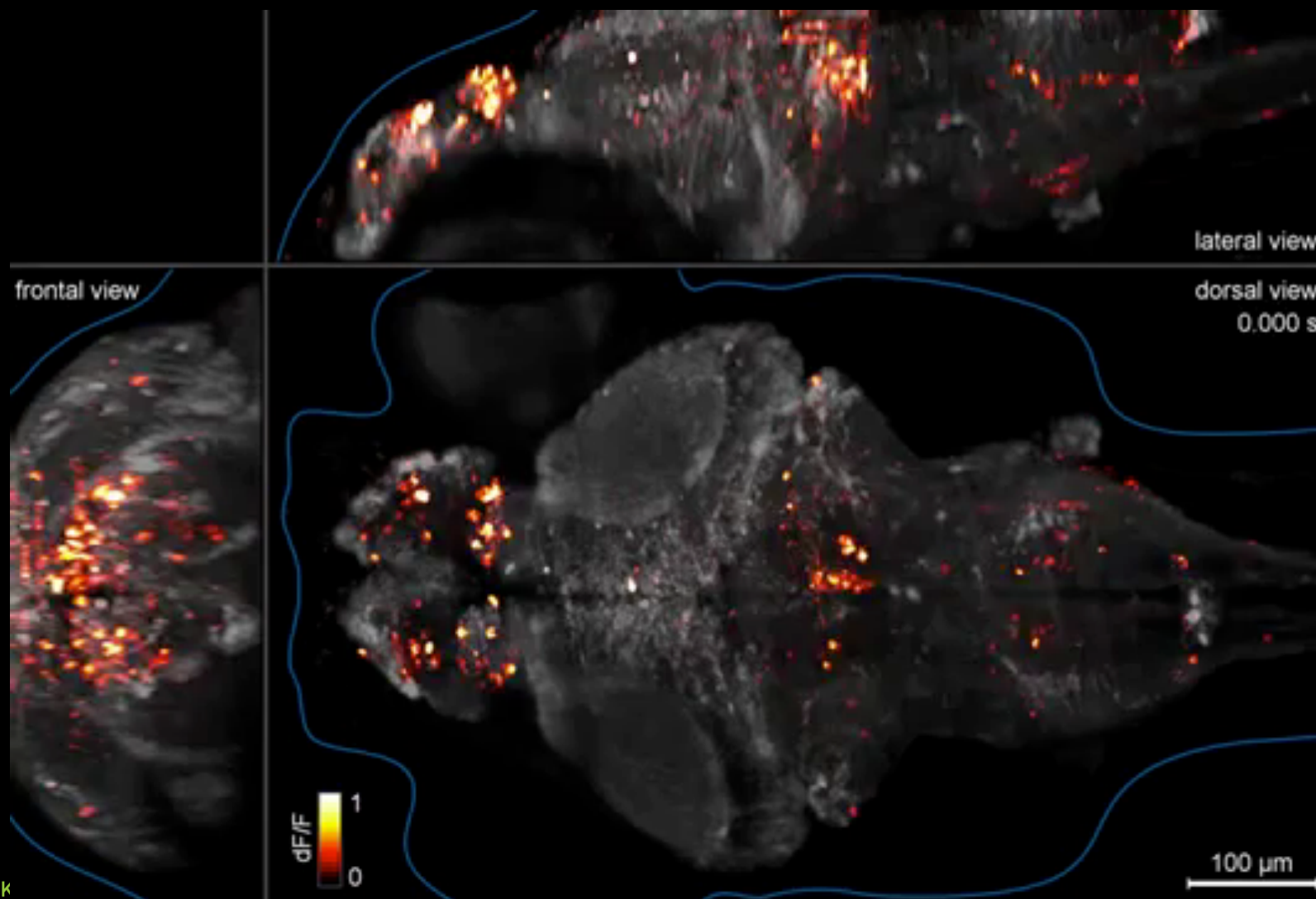
Freeman Lab



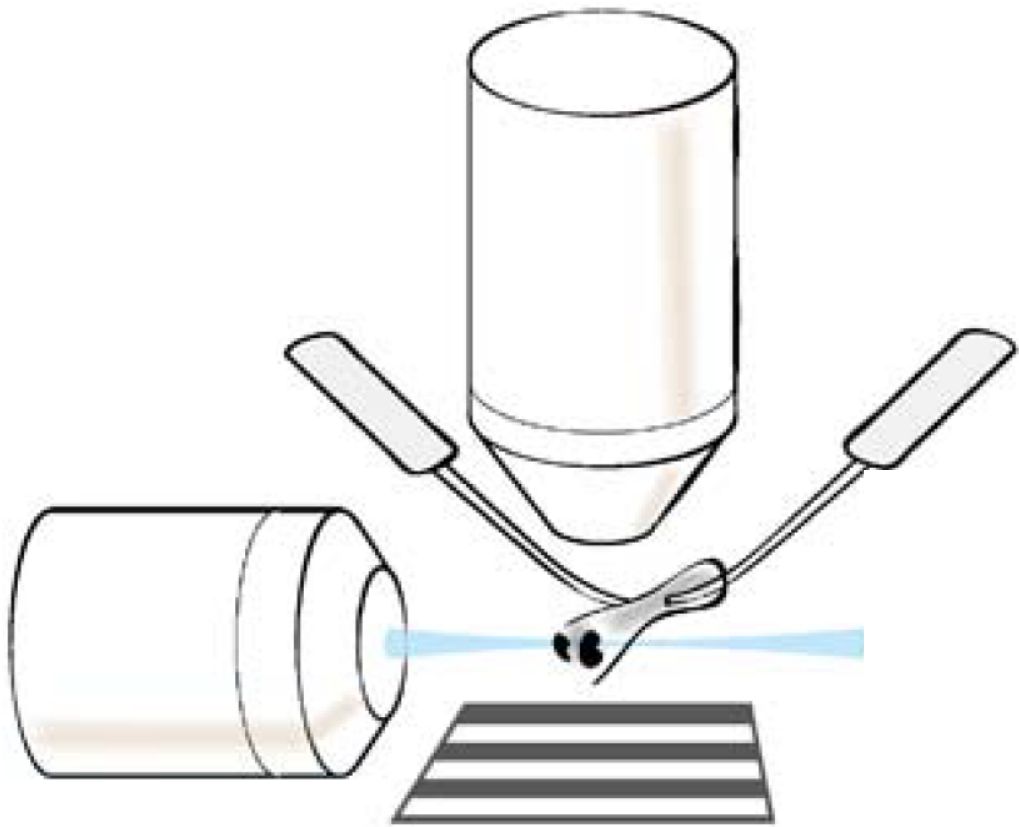
SPARK SUMMIT 2016
DATA SCIENCE AND ENGINEERING AT SCALE
JUNE 6-8, 2016 SAN FRANCISCO



SPARK SUMMIT 2016



SPARK



$t, (x, y, z)$

time $\sim 10^4$

space $\sim 10^7$

$\sim 10^{11}$ elements

~ 1 TB

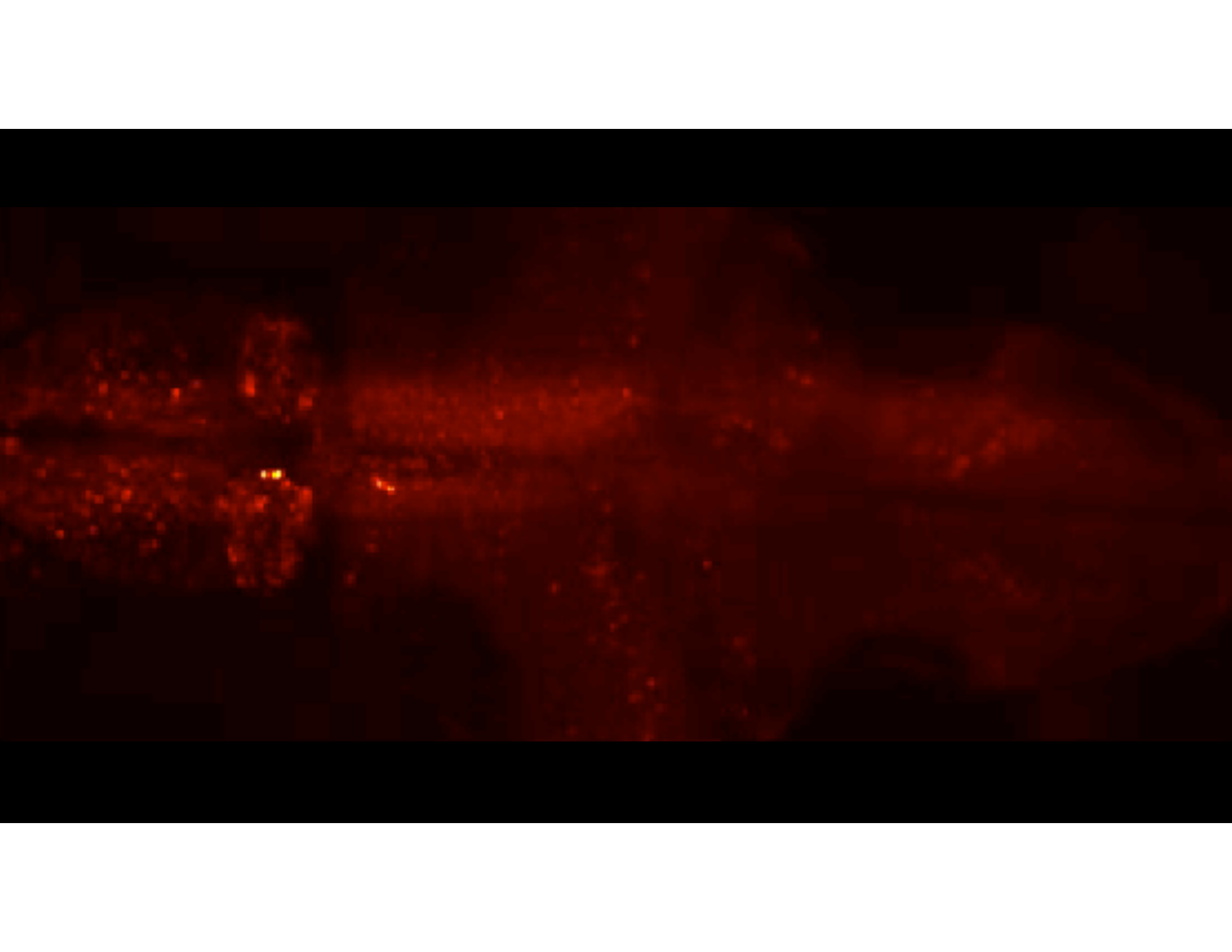
(n, k)

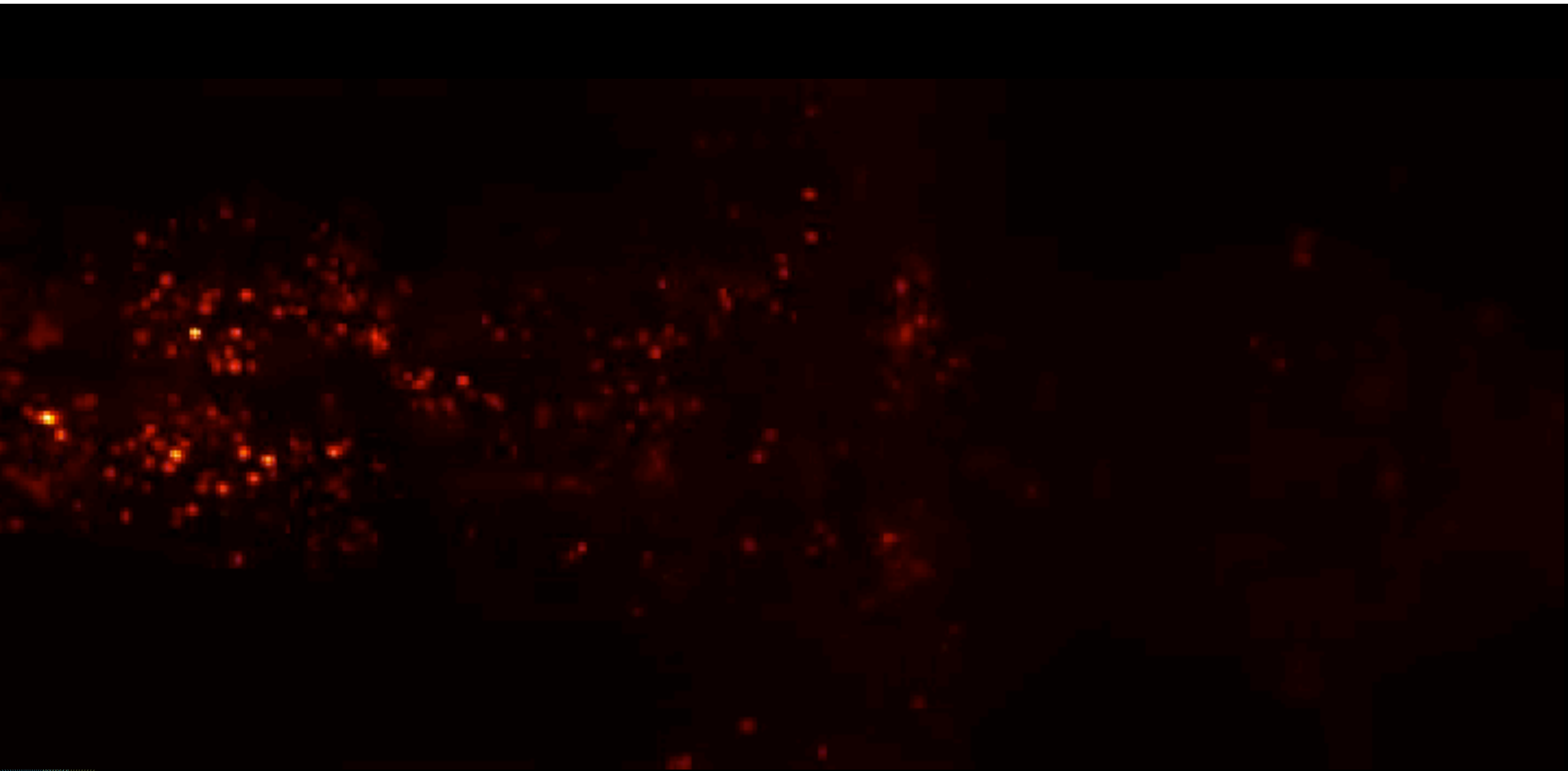
(x, y, t)

(x, y, z, t)

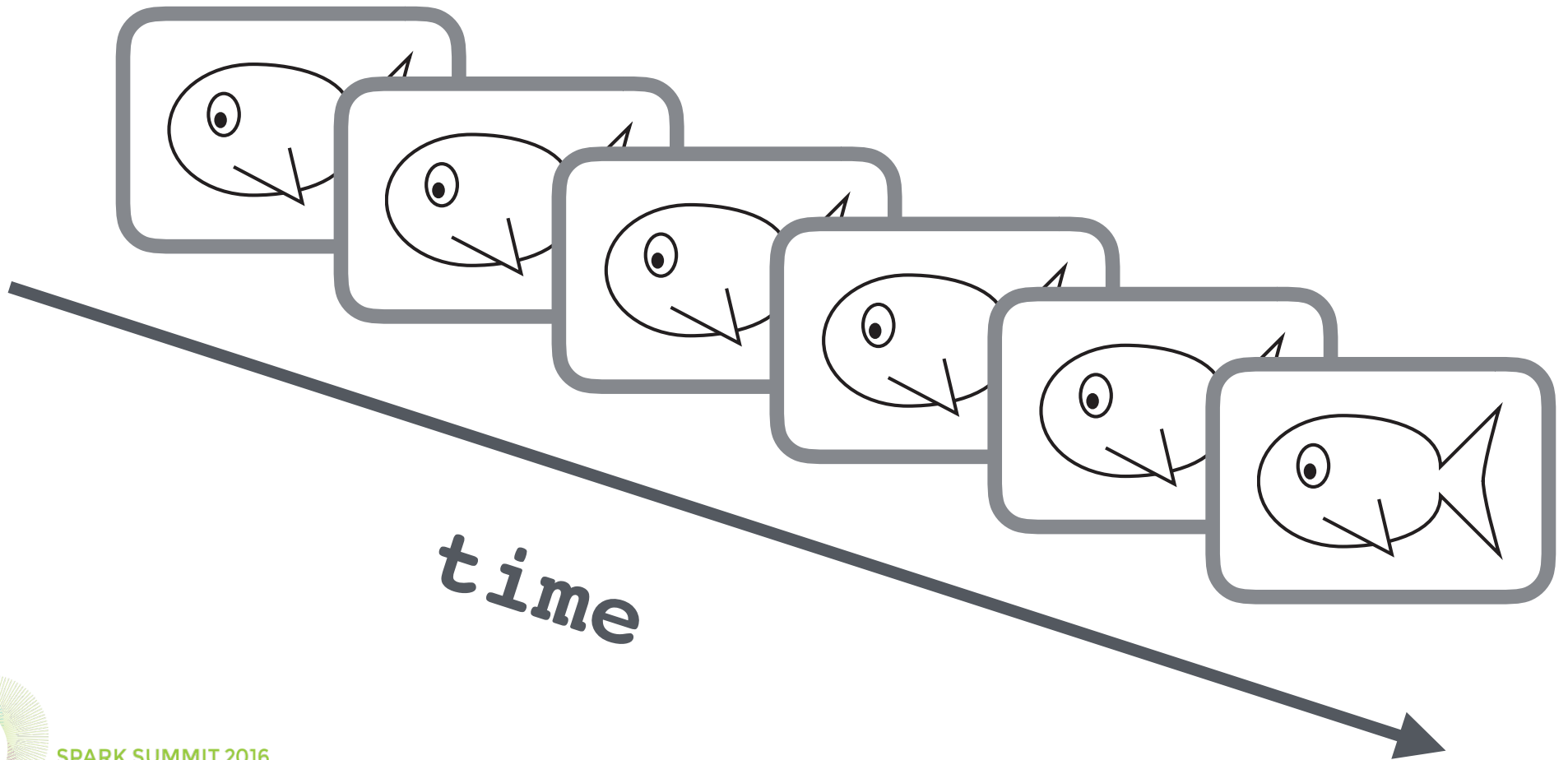
(x, y, z, c, t)



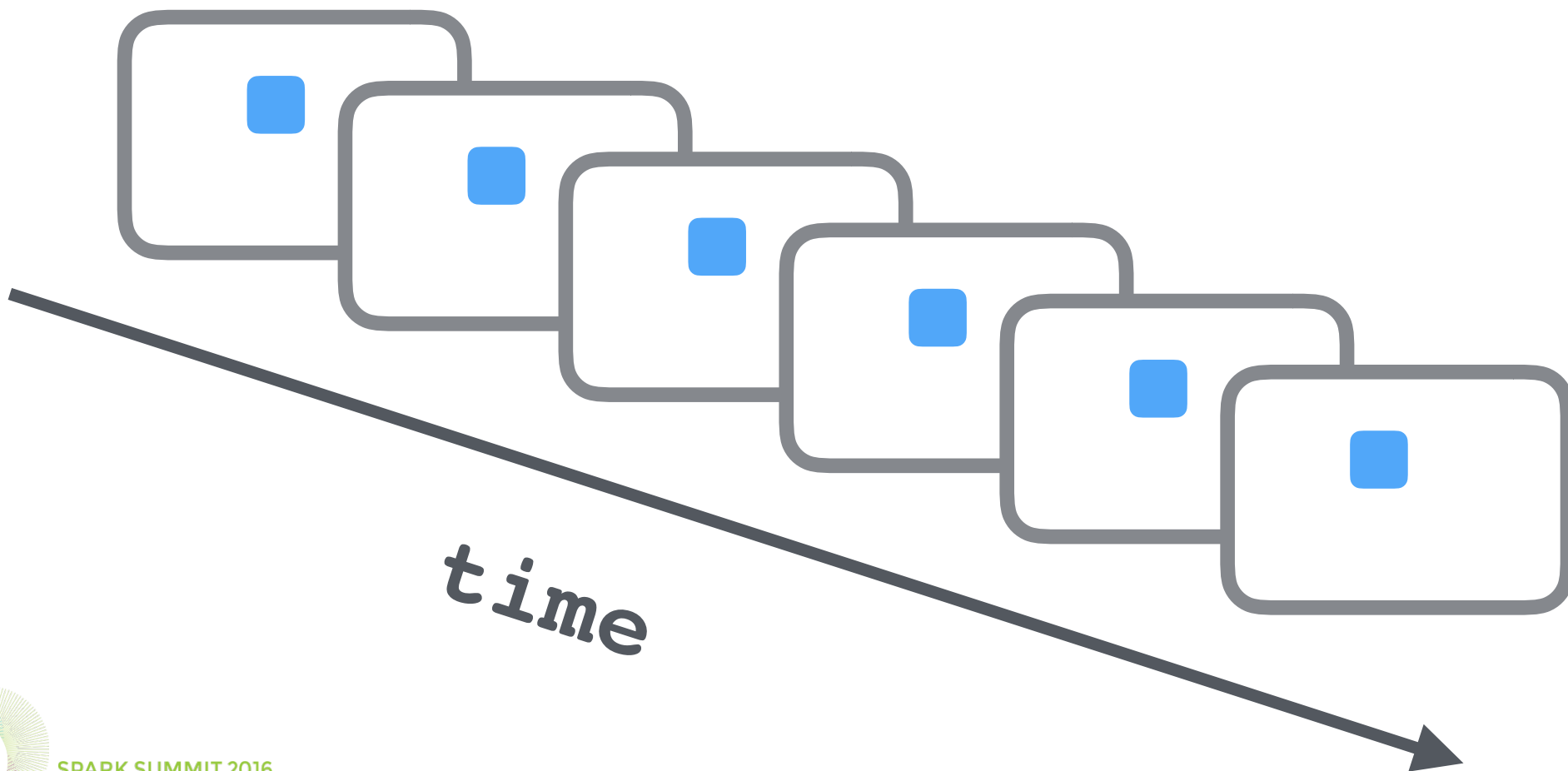




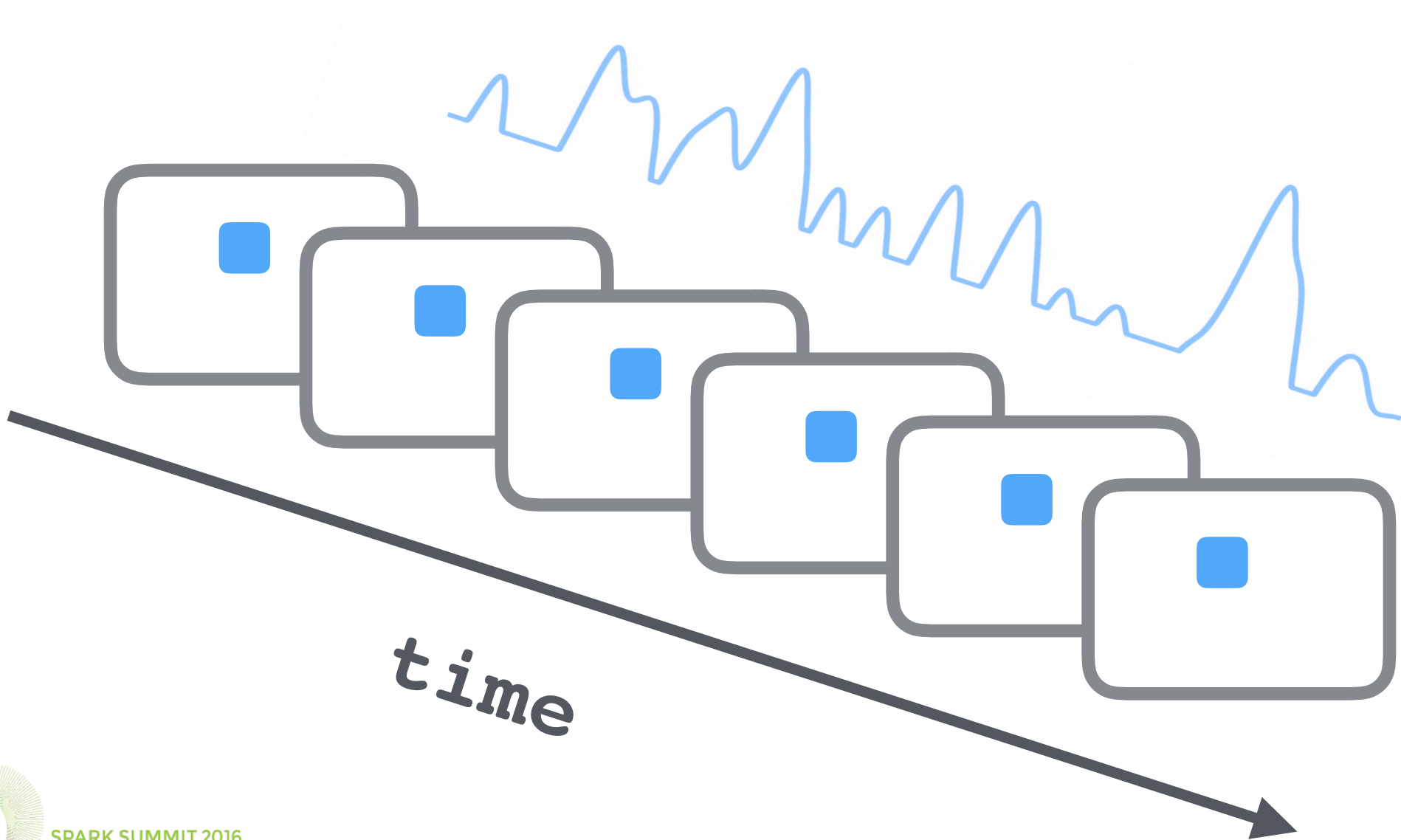
SPARK SUMMIT 2016



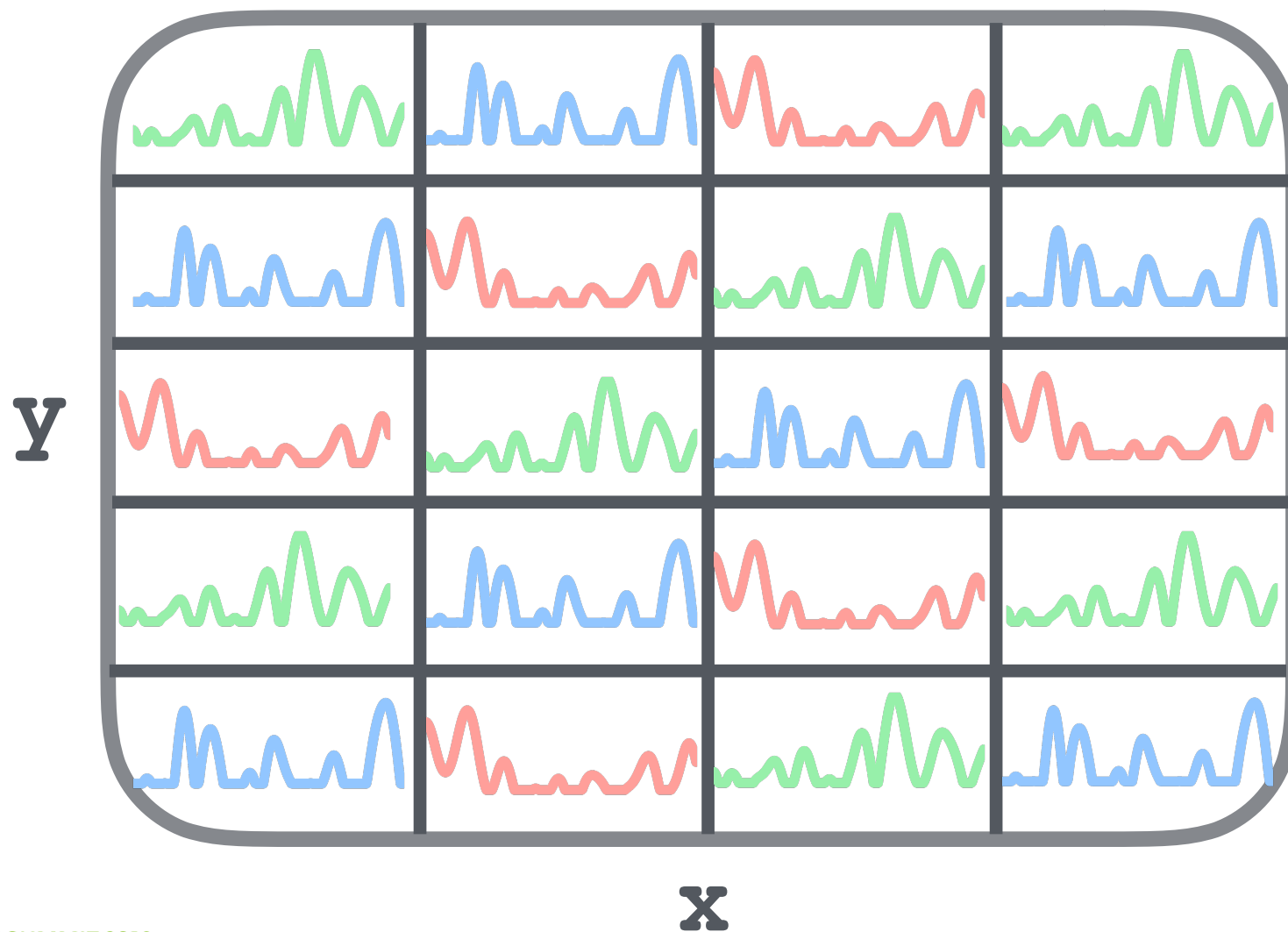
SPARK SUMMIT 2016

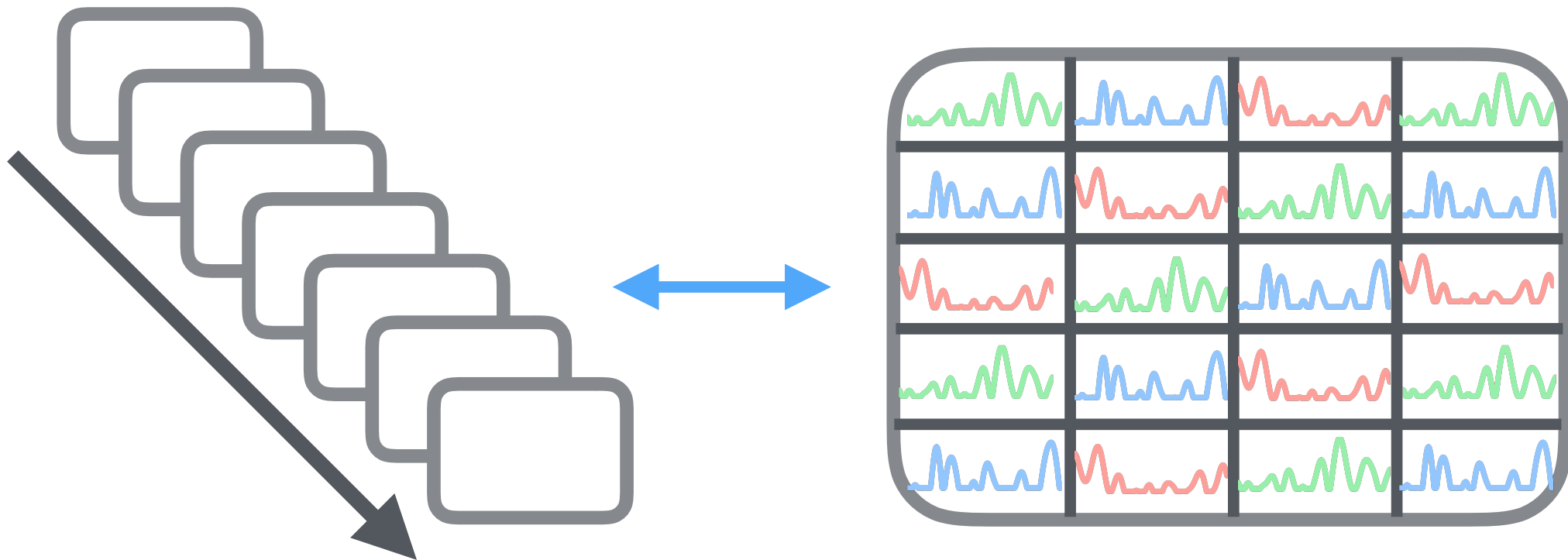


SPARK SUMMIT 2016



SPARK SUMMIT 2016





SPARK SUMMIT 2016

neuroscience

astronomy

geospatial

climate science



SPARK SUMMIT 2016

'bolt

- a distributed ndarray
- built on PySpark
- conforms to NumPy API

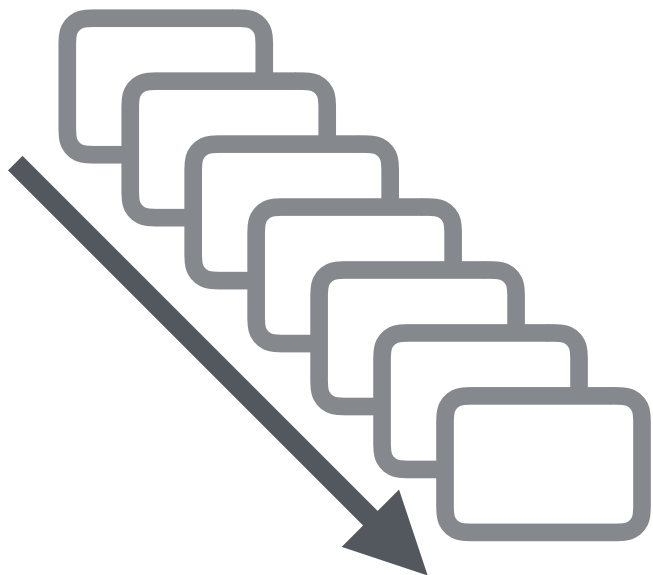


```
data.mean(axis=0)
```

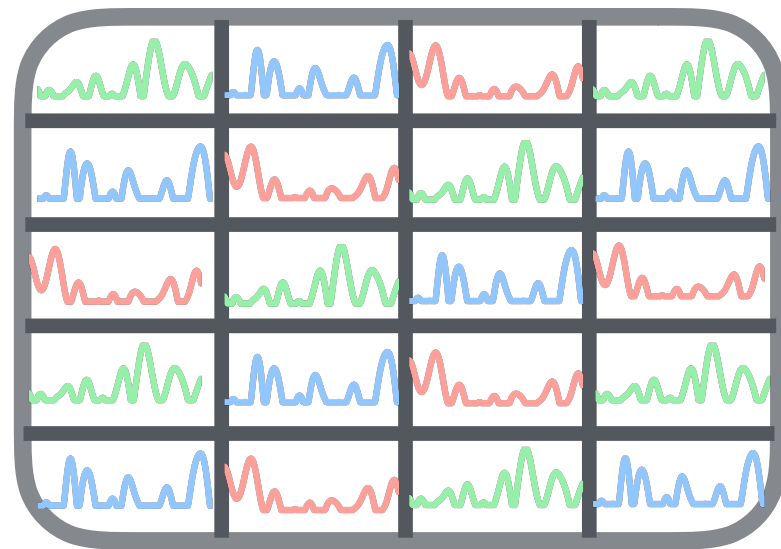
```
data.T
```

```
data[2, 4:10]
```





$(t, | x, y, z)$



$(x, y, z, | t)$

$(v, w, | x, y, z)$

(v, w)

y



x

z

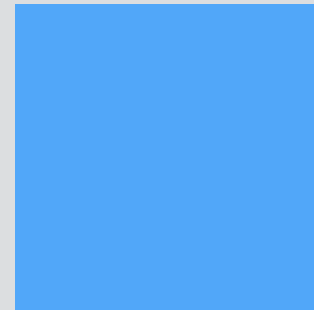


$(v, w, x, | y, z)$

(v, w, x)

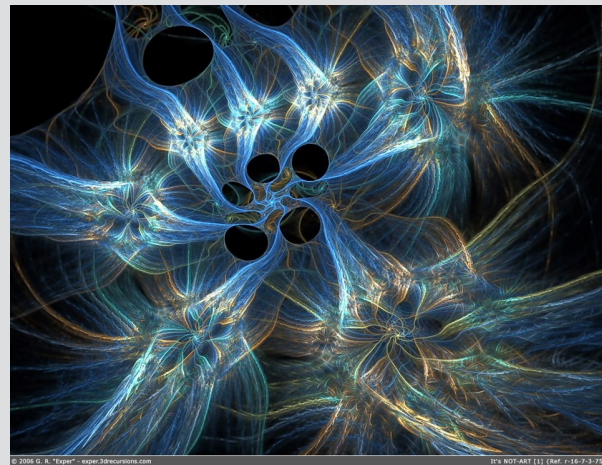
y

z



$(v, | w, x, y, z)$

(v)



SPARK SUMMIT 2016

indexing slicing
apply-along-axis



indexing **slicing**
apply-along-axis
transpose **reshape**



indexing **slicing**
apply-along-axis

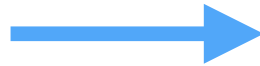
transpose **reshape**

map **reduce** **filter**

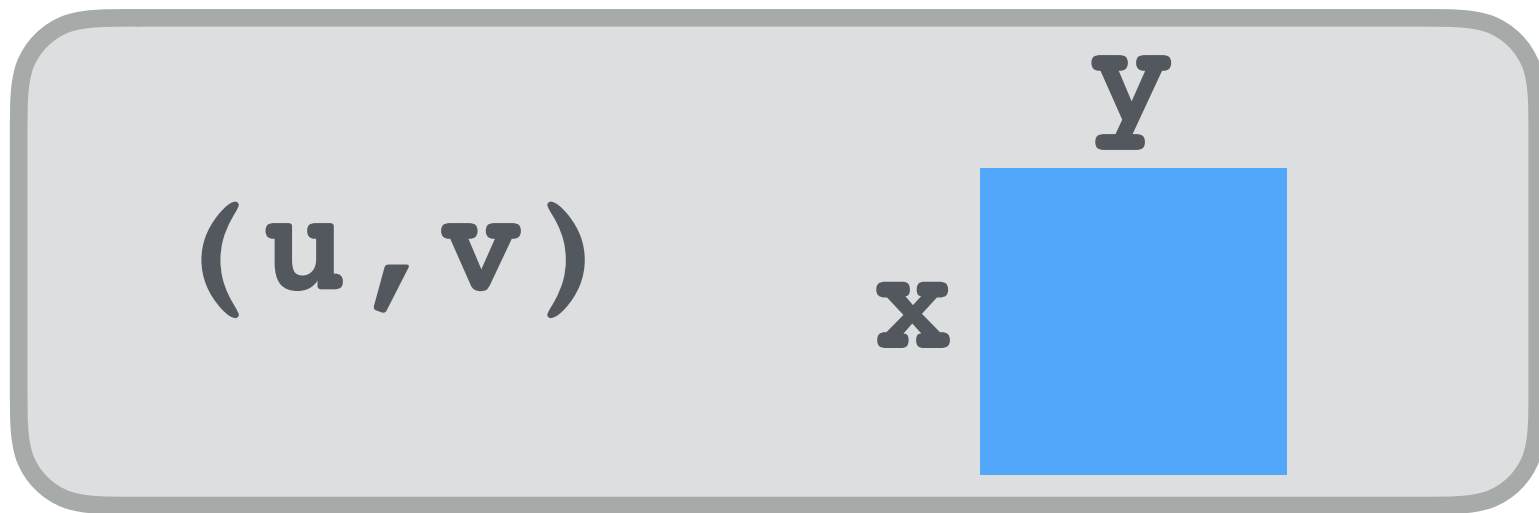
chunking **padding**



indexing



filter
map

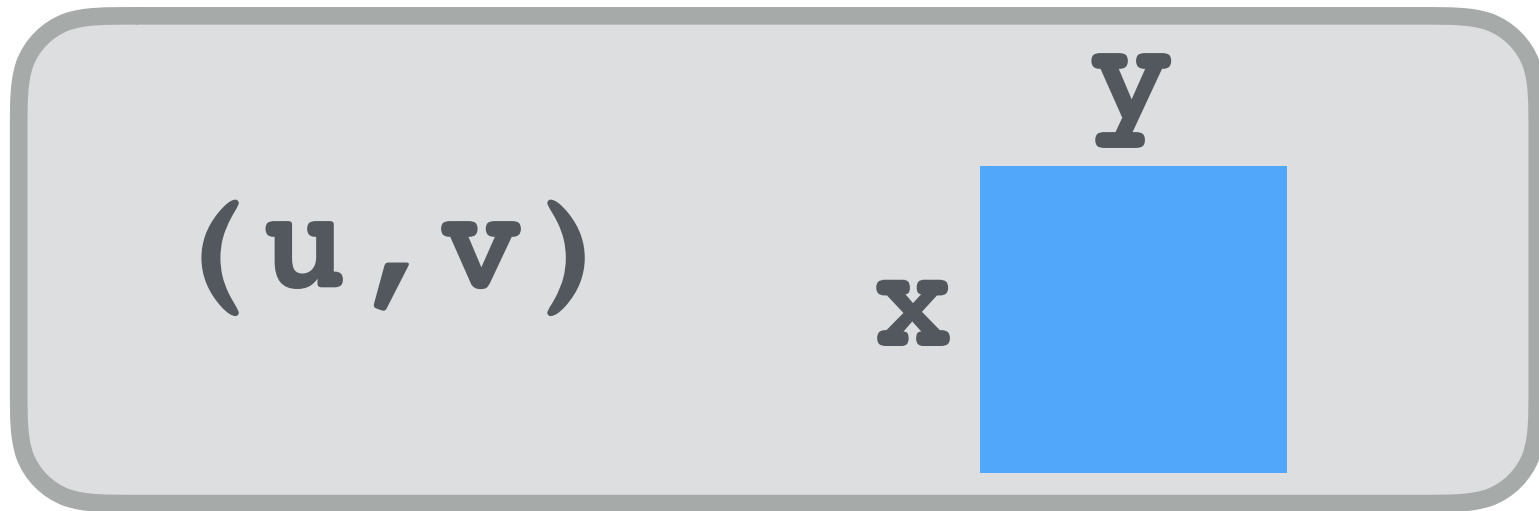


SPARK SUMMIT 2016

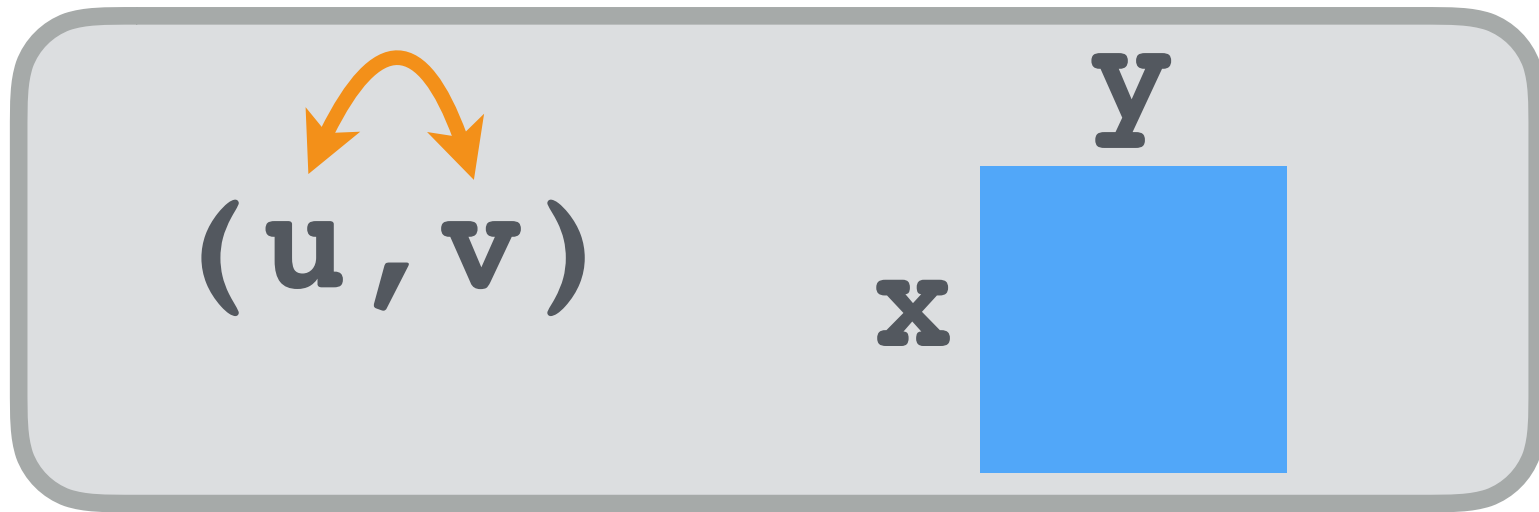
apply-along-axis



map reduceByKey

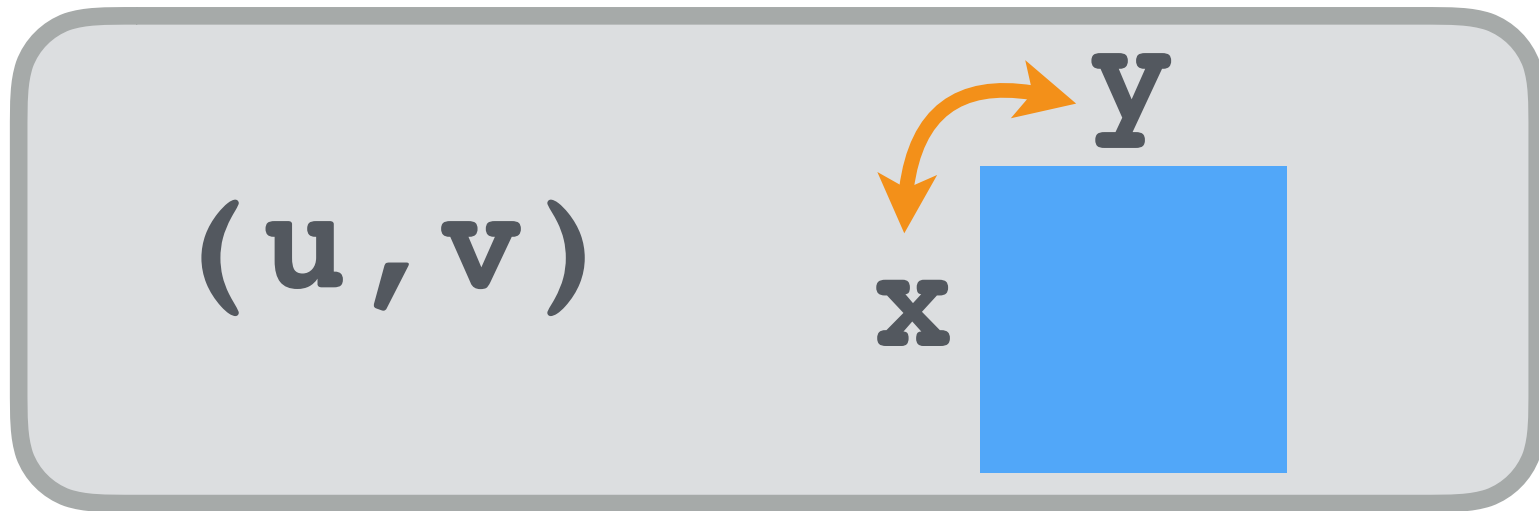


transpose \rightarrow map

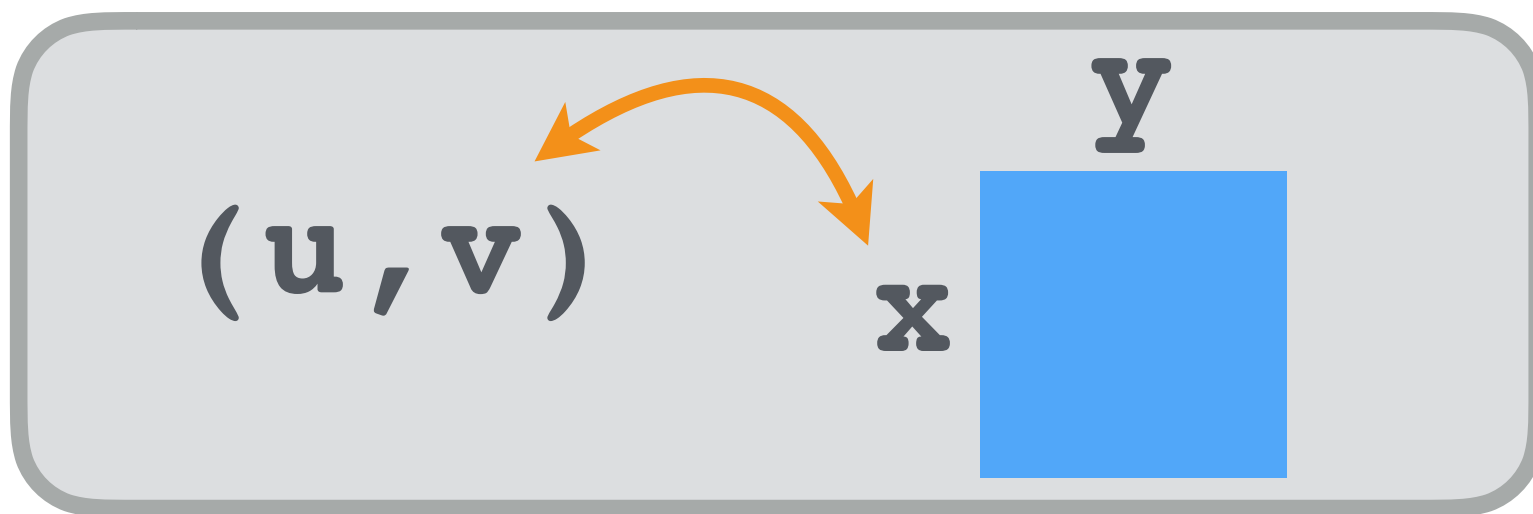


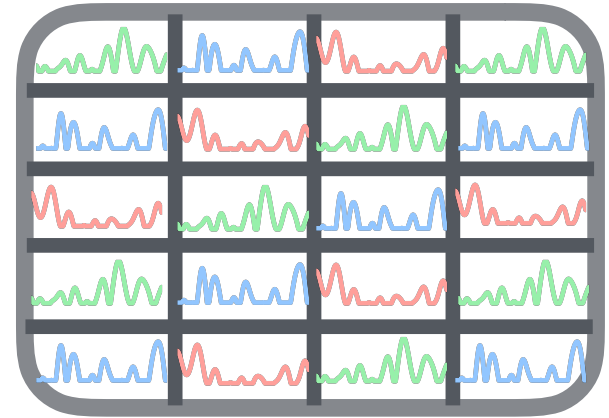
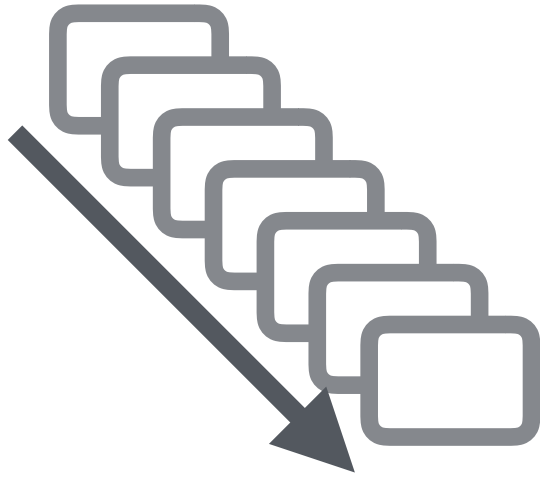
SPARK SUMMIT 2016

transpose → map



transpose → shuffle





$(t \mid x, y, z)$

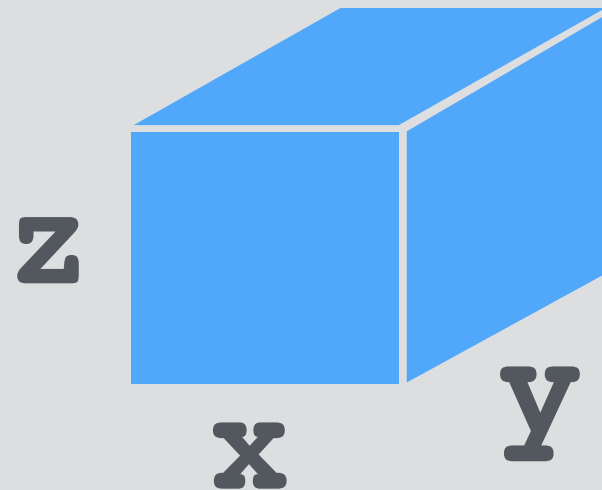
$(x, y, z \mid t)$



SPARK SUMMIT 2016

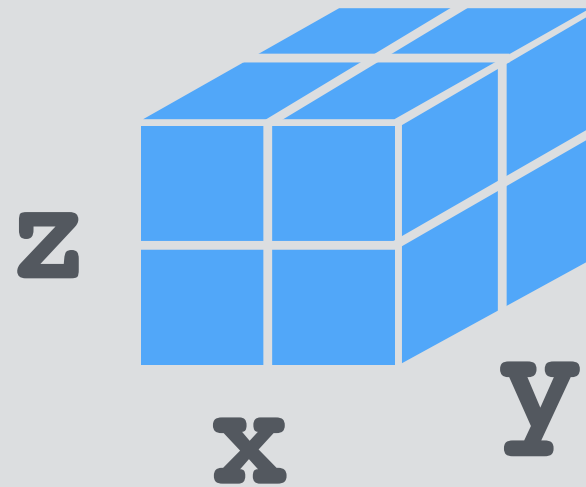
$(t \mid x, y, z)$

(t)



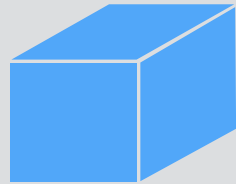
$(t \mid x, y, z)$

(t)

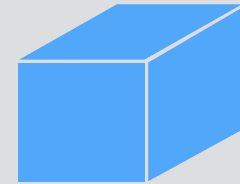


$(t \mid x, y, z)$

(t, chunk)



(t, chunk)



(t, chunk)

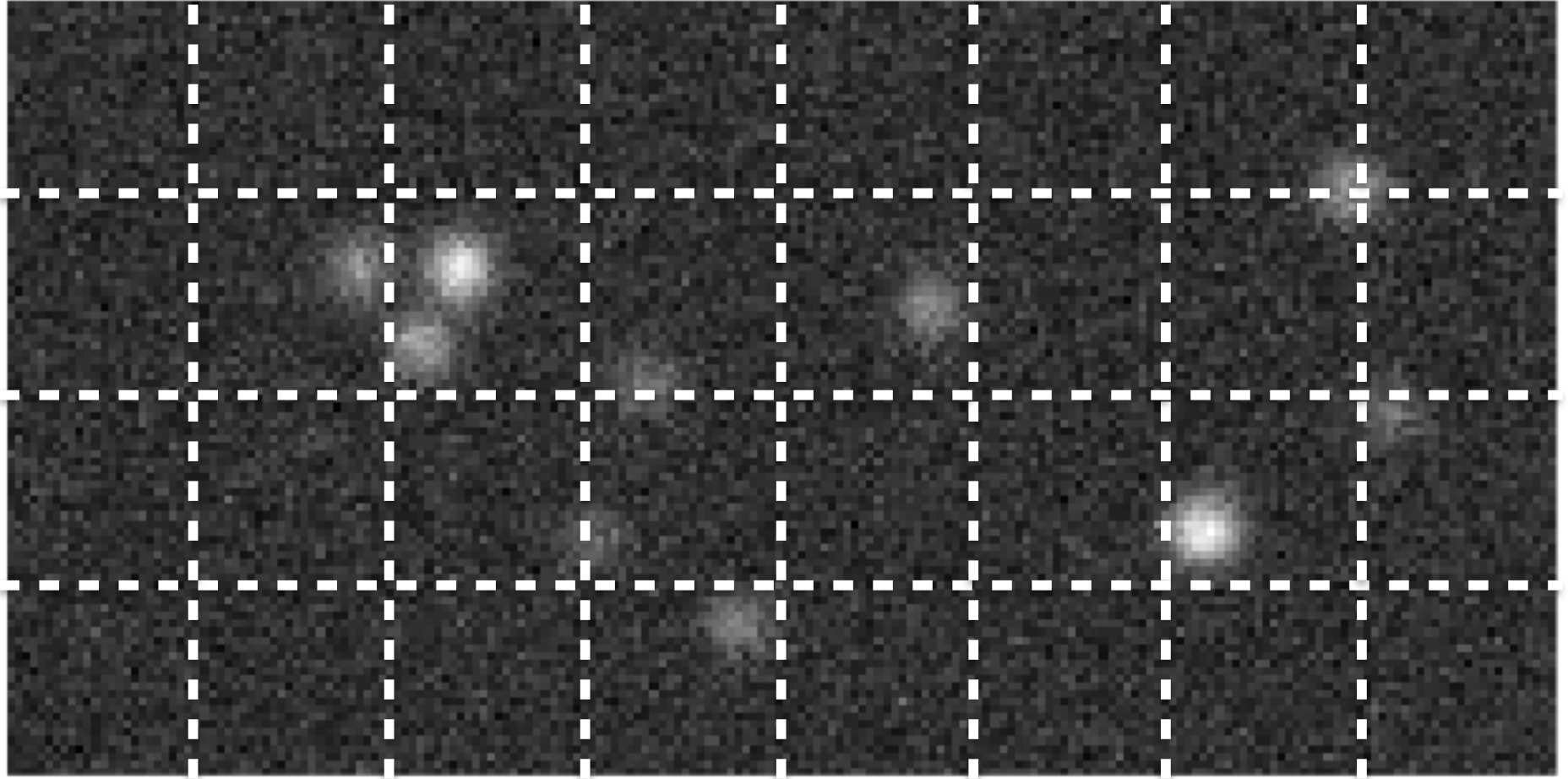


...

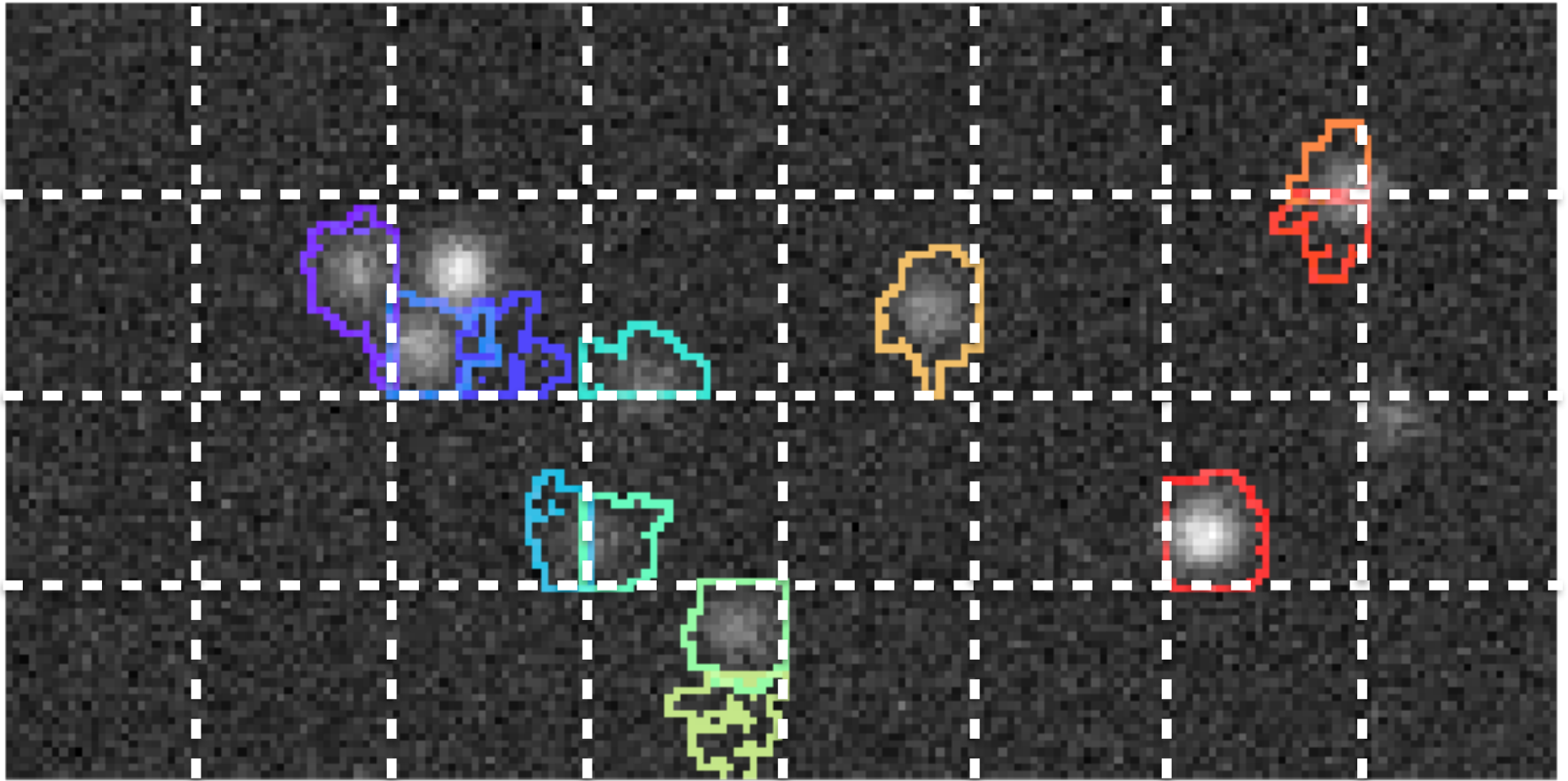




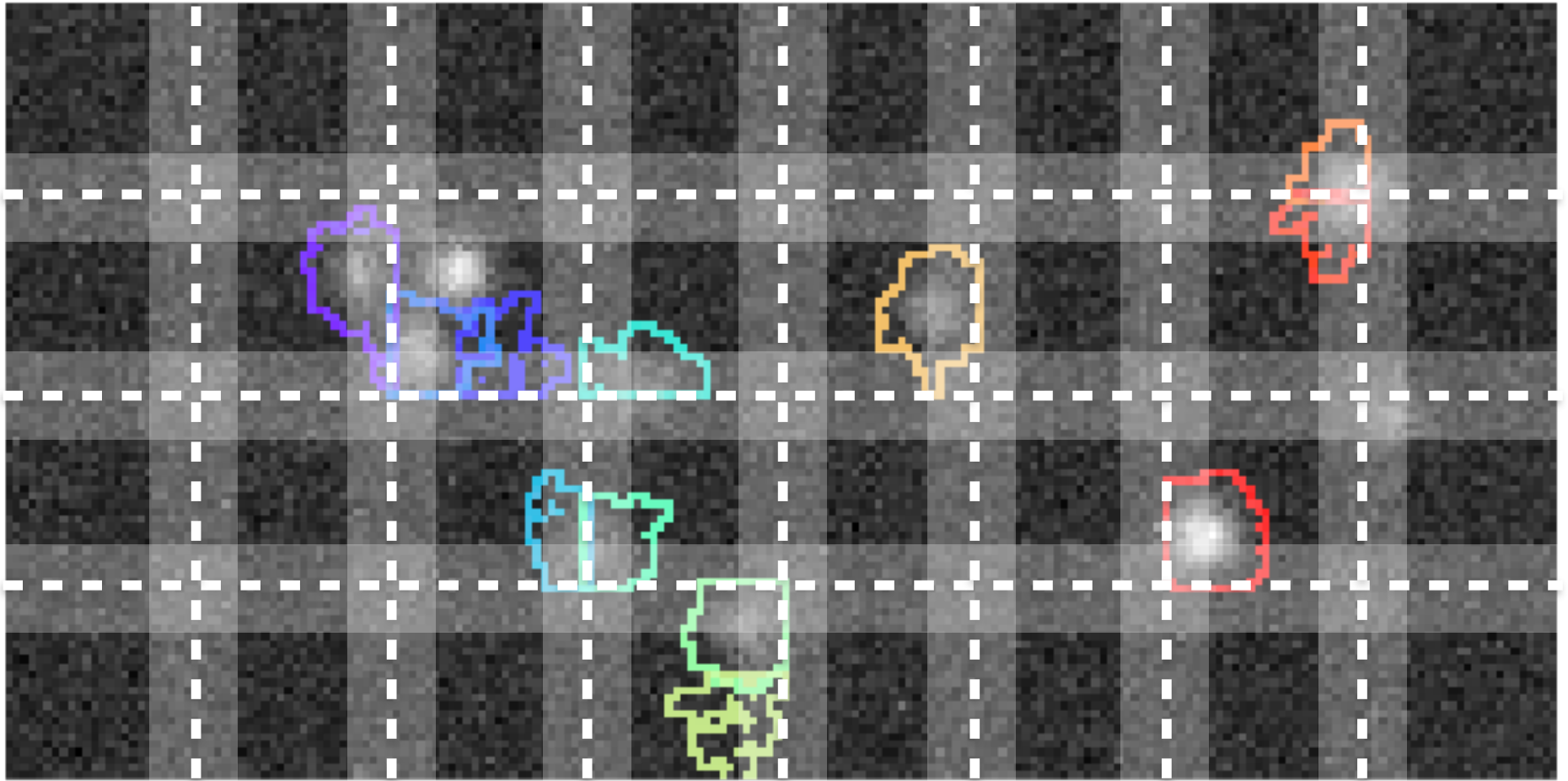
SPARK SUMMIT 2016



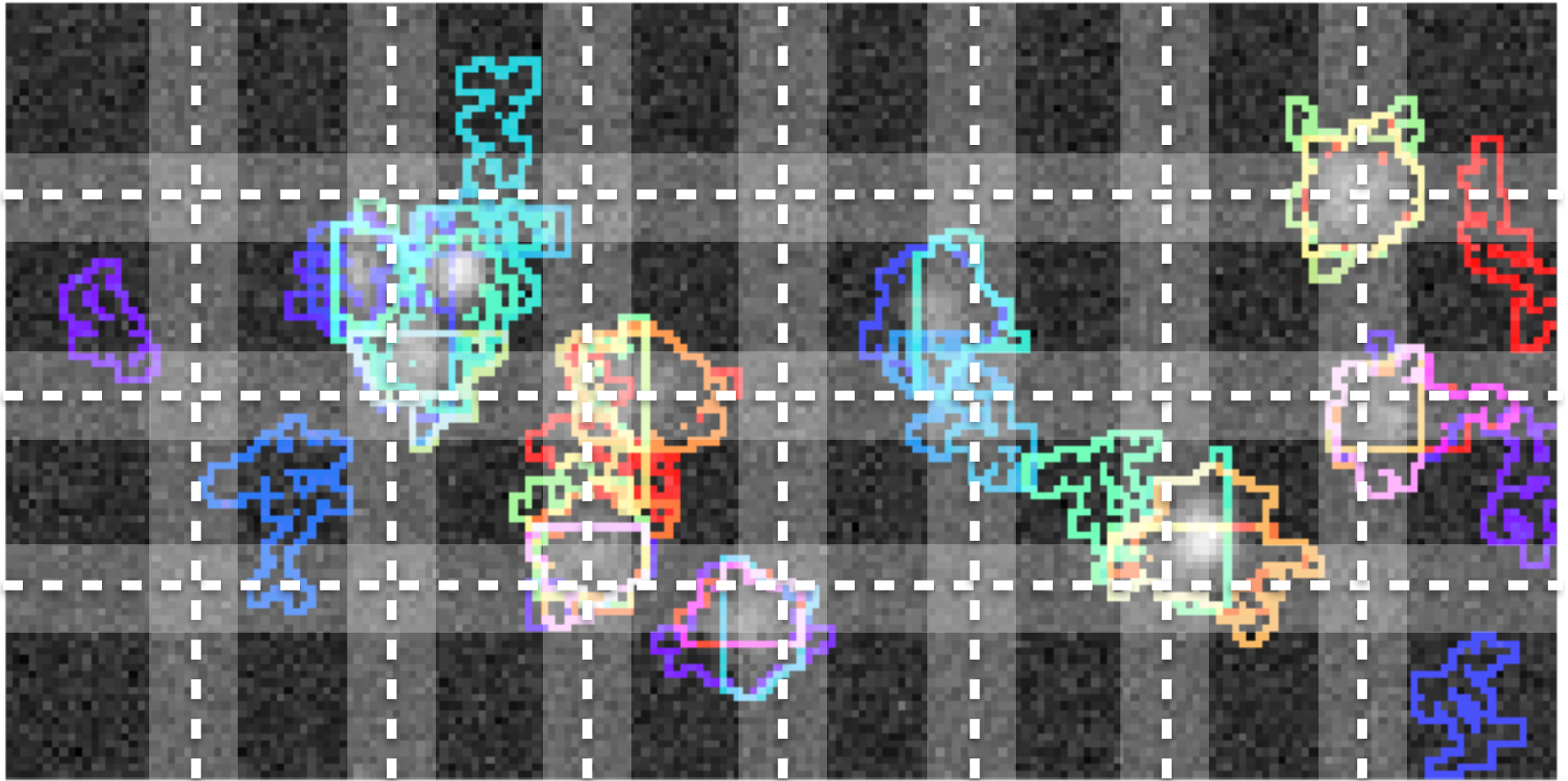
SPARK SUMMIT 2016



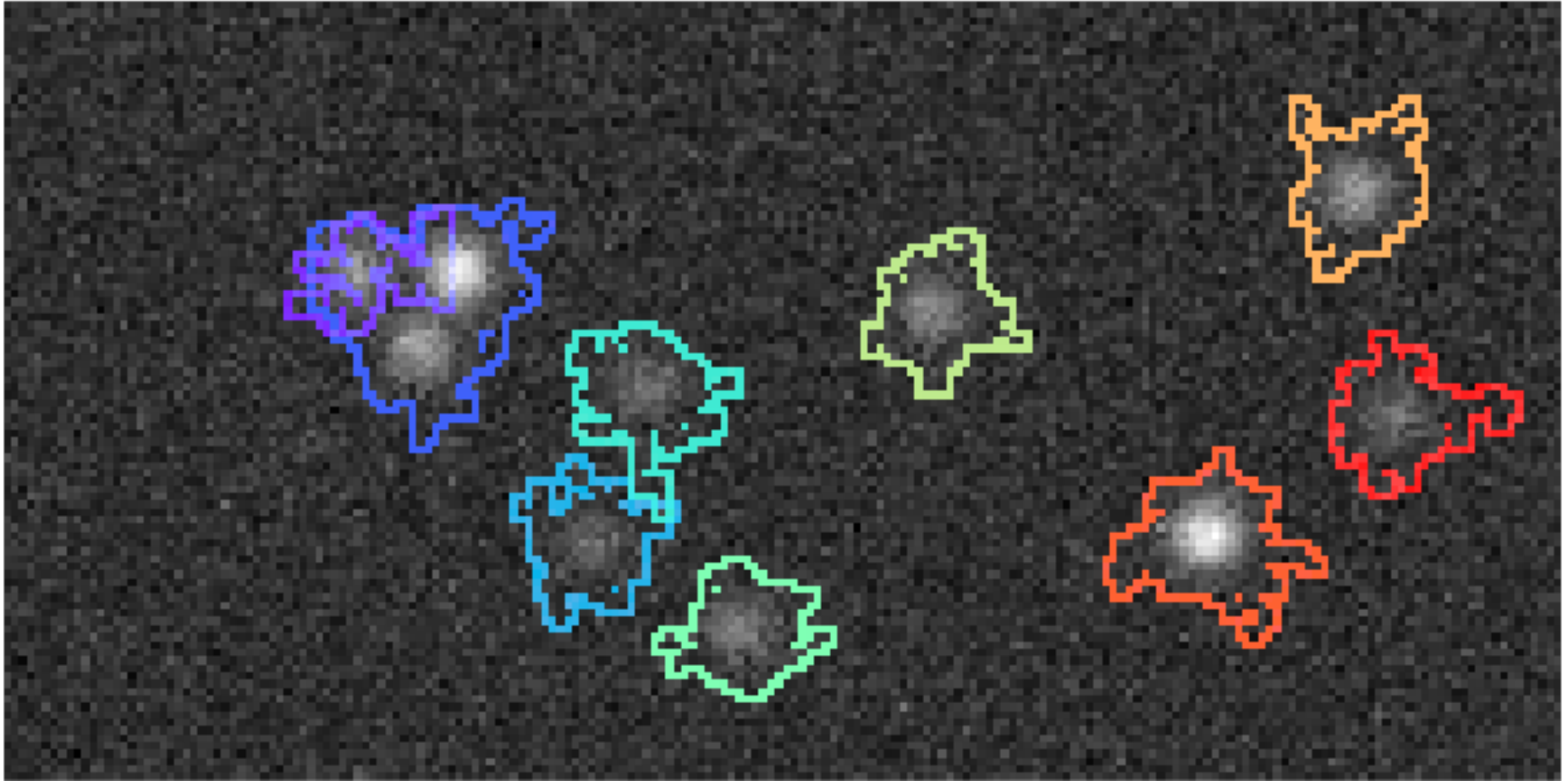
SPARK SUMMIT 2016



SPARK SUMMIT 2016



SPARK SUMMIT 2016



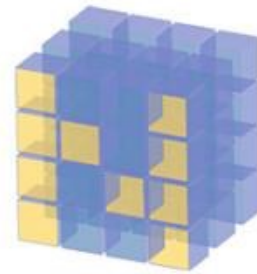
SPARK SUMMIT 2016

chunking
+
transpose = **shuffle
optimization**

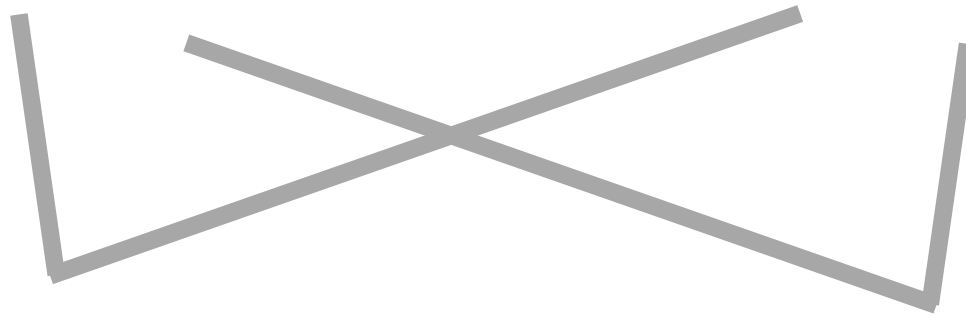


SPARK SUMMIT 2016

'bolt



NumPy



thunder

???



SPARK SUMMIT 2016

thanks

Freeman Lab

Jeremy Freeman
Nicholas Sofroniew
Andrew Osheroff

Janelia Scientific Computing

Ken Carlile
Robert Lines

join us!

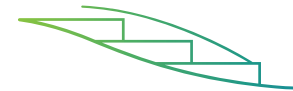
GitHub

bolt-project
thunder-project

Twitter

@jsonWittenbach

hhmi



janelia

Research Campus



SPARK SUMMIT 2016