

Solar Data Analysis Center

I never never use a Dig, Dig D.
— *The First Lord's* song, W.S. Gilbert, HMS
Pinafore



Virtual Solar

What Passes for Big Data in Solar Physics at NASA Goddard

Joseph B. Gurman

NASA Goddard Space Flight Center

Heliophysics Division

Laboratory for Solar Physics

—

Facility Scientist, Solar Data Analysis Center

US Project Scientist for SOHO

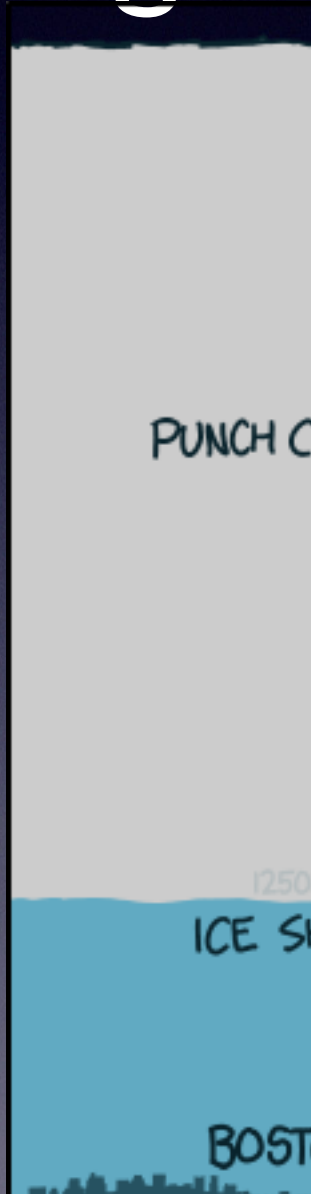
STEREO Project Scientist

What Exactly Do We Mean by “Big?”

Big like Google, whose data holdings have been estimated at ~ 15 Ebyte?

- Or put another way, if encoded on Hollerith cards, one New England 4.5 km deep

Even the NSA probably still holds less.



What Exactly Do We Mean by “Big?” (II)

- The “classic” (2001) definition, which requires:
 - volume
 - variety
 - “velocity” (flux)



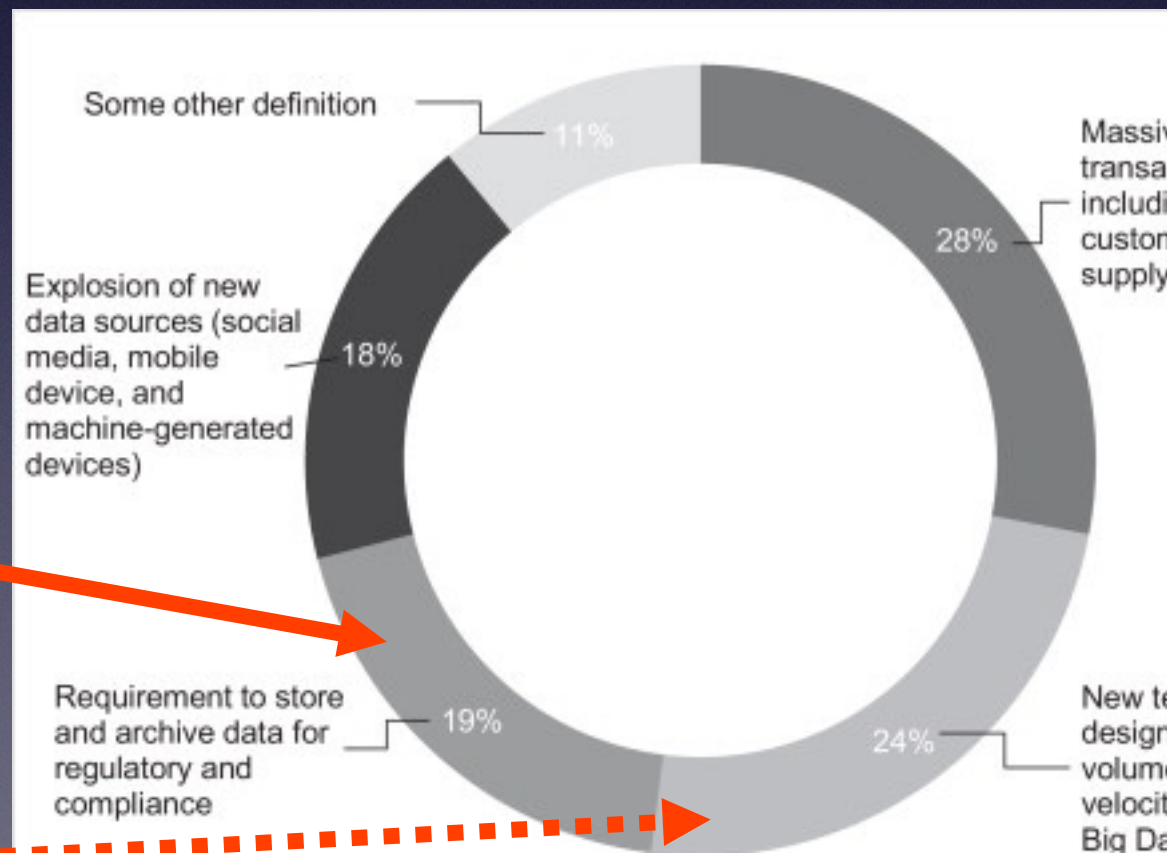
What Exactly Do We Mean by “Big?” (III)

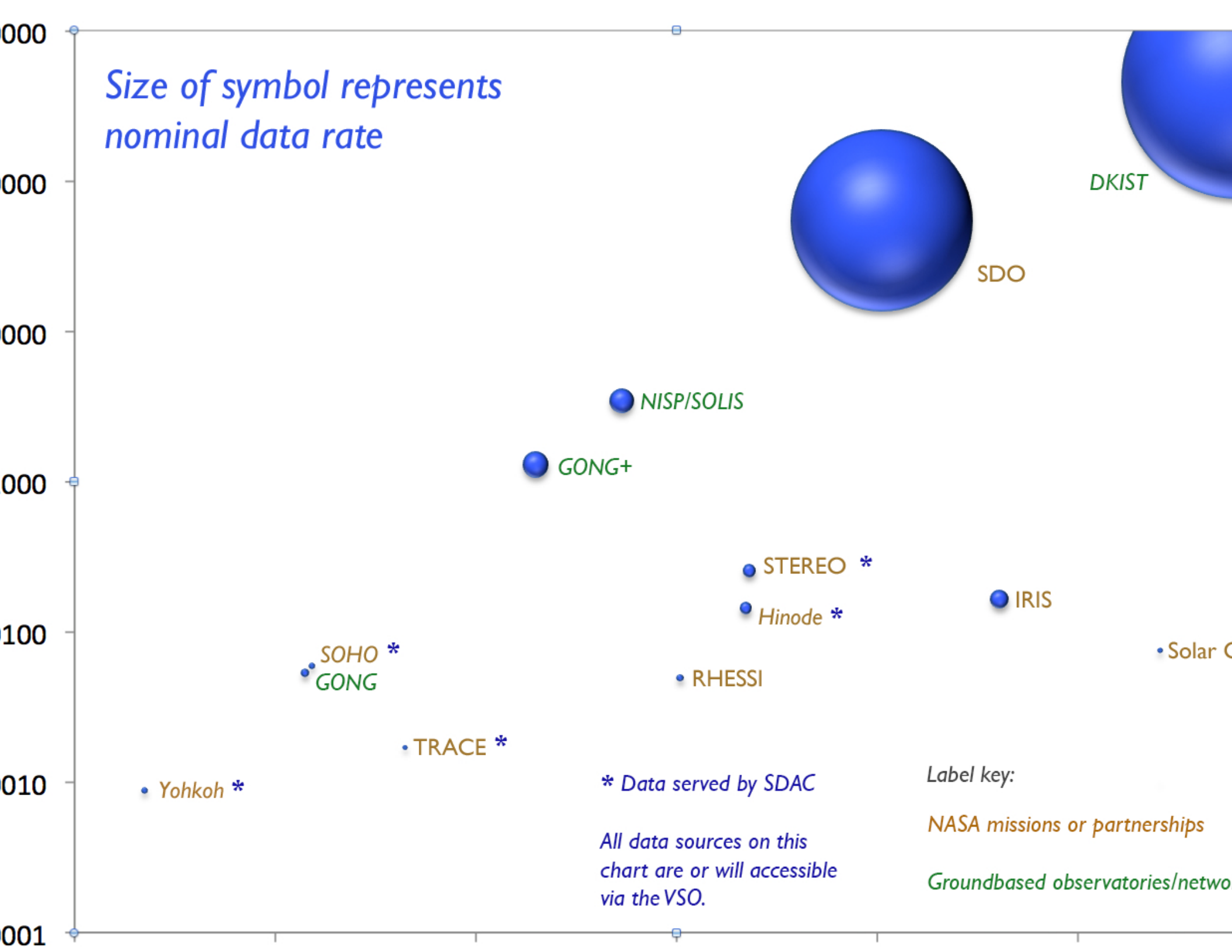
... do we mean “Big Data Methods,” or something else?

... what happens when I ask 154 “C-suite” executives:

... what I’ve been talking about....

... to some extent,





– Leo Byzantius, quoted by Plutarch, Political Precepts

sis Center

Virtual S

We (*the SDAC*) Don't Yet Hold an Entire “Big Data” Archive

Current largest complete archive holdings are H
20 Tbyte) and STEREO (80 Tbyte): small potatoes

largest single-mission holding is SDO (mostly A
ata from periods of interest + last ~ year): 1 Pb

The archive as a whole is characterized by all th
’s (volume, variety, and velocity), but the highe
“velocity” (flux) data set has no variety (SDO A

– Shakespeare, Othello, III, 3

A Story of Days to Come (I)

At some point, entirely TBD, the SDAC is likely to be named the long-term archive for the entire S data set

2 Pbyte/year, 6 years so far

Not supportable by current storage architecture (or room floor)

Comes at a time of planning for change in NASA and data architecture (agency, center, commercial, ? cloud “consolidation”)

– Likely configuration TBD until sometime in FY18

And whistles in his sound.”

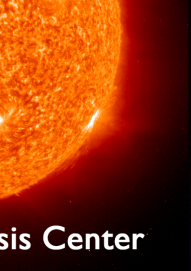
– Shakespeare, As You Like It, II, 7

A Story of Days to Come (II)

wish I could tell you how NASA /
Goddard / the SDAC will be
handling the data from SDO (the
big data solar mission?) three
years from now

Unfortunately, it's unclear and I
have a limited horizon

Somebody else will be making those



Responses to Questions (I)

the processes for planning for future (5-10 years) capabilities of your service? How do you gather input for this planning process and where does input typically come from? Which new features have highest priority?

CS is mission-oriented; meets with mission scientists to understand service; VSO is user-oriented, using frequent outreach and frequent meetings to identify data provides and services to a

CS new features are typically driven by the available technology; VSO features represent what we develop (or add) in response to user community requests.

– David Packard

Virtual S

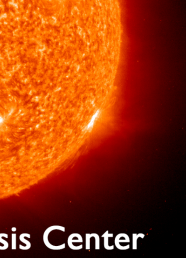
Responses to Questions (II)

ture(s) of your service would you like to stop performing? How do you gather in-
ch decisions and where does input typically come from? What is preventing you

would truly, deeply like to get out of the business of
ing new storage on multiple tiers every five to seven y

it” is experience-based, by the facility scientist. Wasteful
ue battling numerous, partially conflicting agency direct

of clearly effective alternatives that meet mission risk
ement criteria and a current period of confusion (“des



sis Center

– David Packard



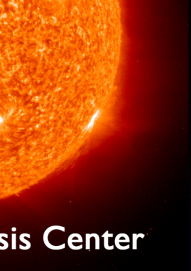
Virtual S

Responses to Questions (III)

steps you are taking to make your data interoperable with allied data sets from other agencies and out of NASA? How do you find allied data sets and what criteria make data sets suitable for enabling interoperability?

Our solar data pretty much are currently interoperable with any other solar physics data (FITS format; software libraries in Python). Interoperability with other scientific disciplines, particularly heliophysics, requires specialist knowledge, so not clear whether investing in further interoperability buys anything.

Virtual Solar Observatory: all solar physics data, not just space-based
SPDF + VxOs: access to heliophysics data



sis Center

scientific analysis



Virtual S

Even the Virtual Solar Observatory (VSO)'s data dictionary, which recognizes only 16 physical observables, is in a data science sense too rich

Only three basic organizations of solar data: map (image data organized as one), spectrum, and light curve

Real variety is in the solar features analyzed and the methods used to analyze them

If you want to populate the variety dimension of a NASA solar big data map, look at the work of the SDO Feature Finding Team

- distributed effort to develop multiple big data tools
- http://solar.physics.montana.edu/sol_phys/fft/static/