# First results from PARSE.Insight
## The HEP survey on
## data preservation, re-use and (open) access

- **Background information**
- **The HEP Survey**
- **Next steps**

Andre Holzner (CERN), Peter Igo-Kemenes (Gjøvik/CERN), Salvatore Mele (CERN)

First Workshop on Data Preservation and Long Term Analysis in HEP
[Open Access and Long Term Collaborative Governance]
DESY - 26-28 January 2009

# Actors

**Alliance for Permanent Access**

Association of stakeholders with an interest in digital preservation of publications and research data.

- Research institutes (CERN, ESA, STFC, Helmholtz, MPG)
- Publishers (STM association)
- Libraries (DNB, KB, BL)

CERN involvement as natural extension of Open Access mission and exploration of possible opportunities

After connectivity (GEANT) and Grid (EGEE), both seen as parts of a e-infrastructure, and development of "repositories" (=databases) for publications,FP7 might look at data and data preservation as next frontier and need information to make strategic funding decision

# PARSE.Insight

http://www.parse-insight.eu/

- Small (1.2M€) FP7 project (CSA) with partners from science (CERN, ESA, STFC, MPG) corporations (STM) and libraries (KB,DNB,…). 3/2008 to 2/2010. Lead by STFC (CASPAR, OAIS)

- Interdisciplinary study to offer an insight on:
  - "who is doing what in the field of digital preservation [of research data] and why they are doing it that way"
  - "gap" between what should be done and what is done

- All-around approach and case studies. HEP is the largest case study.

- Main deliverables (to inform FP7 policies/strategies):
  - Insight and Roadmap on issues/threats/opportunities in digital preservation
  - (Gap analysis: what's there what should be there)
  - Example from some communities (N.B. HEP innovator elsewhere!)

# The PARSE.Insight HEP Case Study

## What it is not:
- Technical solutions for data formats
- Technical solutions for data migrations
- This workshop!

## What it is:
- Motivations *vs.* Concerns
- Threats *vs.* Opportunities
- Wishes *vs.* Obstacles

## Our target:
- Make the scientific case with FP7 to use HEP as a case study for a e-infrastructure pilot in preservation ?
- Gather evidence on attitudes of the community to be used when and if we will have to make policy decision on access

# The PARSE.Insight HEP Case Study

Three-pronged approach (Sept. 2007)
- Large-scale survey of the HEP community
- Follow-up and ad-hoc interviews to go into
  - technical details
  - approaches of past and current experiments
  - "superusers" vision
- Tripartite workshop with data producers, consumers, IT support to expose "gaps" between wishes, needs and reality

Re-scoping for synergy with this/these workshop(s)
- Large-scale survey untouched.
- On-demand and ad-hoc interviews to deepen and complete the results of the survey
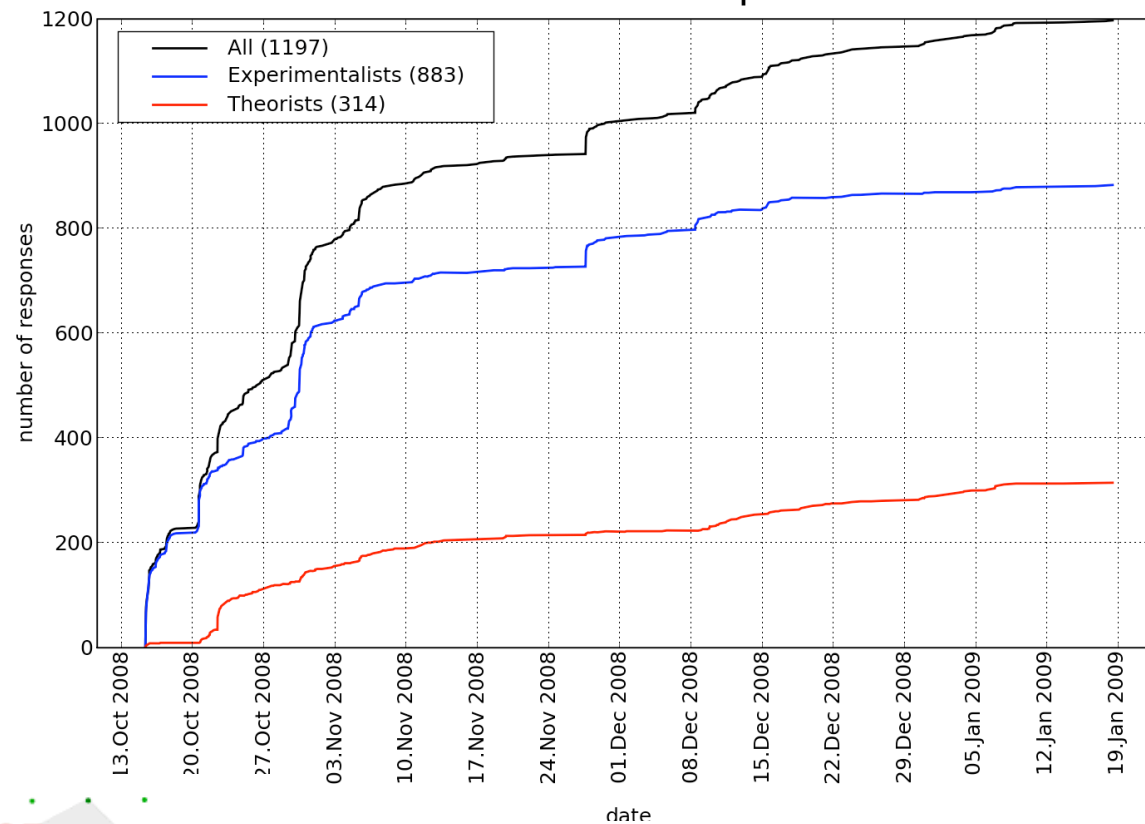- Workshop to be re-scoped following our discussions here

## Today: first results from the survey

# Survey strategy and response

- Livetime of 3 months 10/08-01/09
- Advertised on collaboration and theory mailing lists
- One or two reminders, according to response monitoring
- Advertised on SPIRES twice (reach more spread theorists)
- 1'200 answers (74% exp, 25% th). Target size ~20'000-30'000
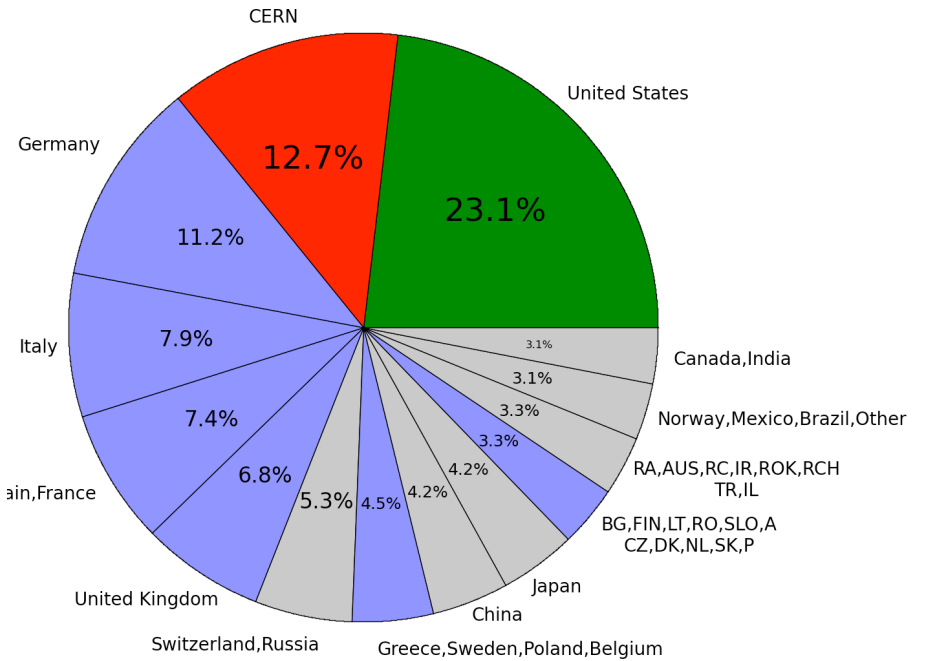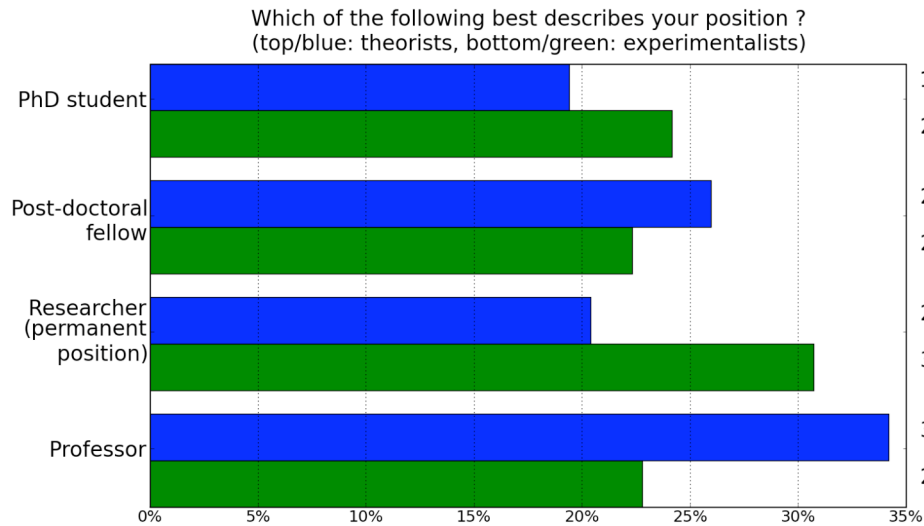


Cumulative number of responses vs. time

# Survey structure

1. Demographics
2. The importance of preservation
3. What to preserve
4. When, how and where to preserve it
5. Threats

# Survey demographics

## Reflects the demographics the community



Which of the following best describes your position ?
(top/blue: theorists, bottom/green: experimentalists)

| Position | Theorist | Experimentalist |
|---|---|---|
| PhD student | 19.4% | 24.2% |
| Post-doctoral fellow | 26.0% | 22.3% |
| Researcher (permanent position) | 20.4% | 30.7% |
| Professor | 34.2% | 22.8% |

CERN 12.7%
United States 23.1%
Germany 11.2%
Italy 7.9%
7.4%
6.8%
5.3% 4.5% 4.2% 4.2%
3.3%
3.3%
3.1%
3.1%
Canada,India
Norway,Mexico,Brazil,Other
RA,AUS,RC,IR,ROK,RCH TR,IL
BG,FIN,LT,RO,SLO,A CZ,DK,NL,SK,P
Japan
China
Greece,Sweden,Poland,Belgium
United Kingdom
Switzerland,Russia
Spain,France
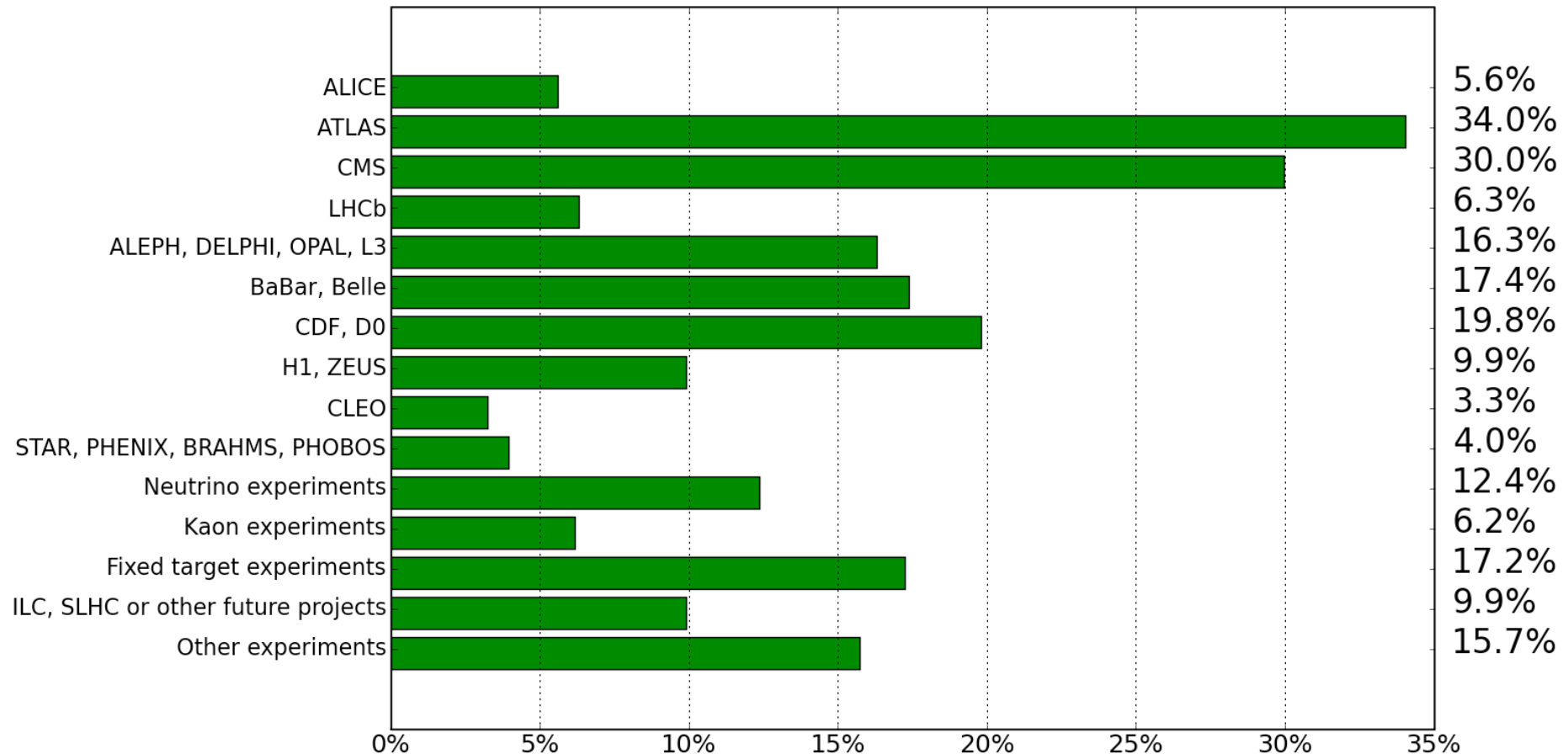
# Survey demographics

## Reflects the demographics of experiments

In which experiments are you / have you been involved ?



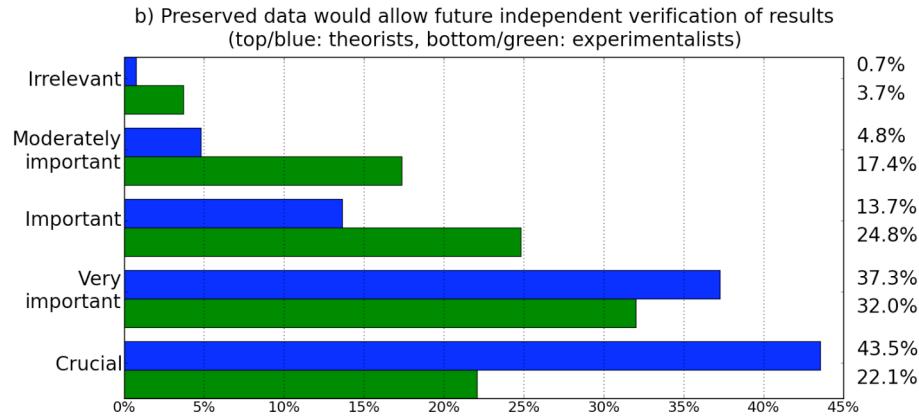| Experiment | Percentage |
|---|---|
| ALICE | 5.6% |
| ATLAS | 34.0% |
| CMS | 30.0% |
| LHCb | 6.3% |
| ALEPH, DELPHI, OPAL, L3 | 16.3% |
| BaBar, Belle | 17.4% |
| CDF, D0 | 19.8% |
| H1, ZEUS | 9.9% |
| CLEO | 3.3% |
| STAR, PHENIX, BRAHMS, PHOBOS | 4.0% |
| Neutrino experiments | 12.4% |
| Kaon experiments | 6.2% |
| Fixed target experiments | 17.2% |
| ILC, SLHC or other future projects | 9.9% |
| Other experiments | 15.7% |

# The importance of preservation

In your opinion, how important is the issue of data preservation ?
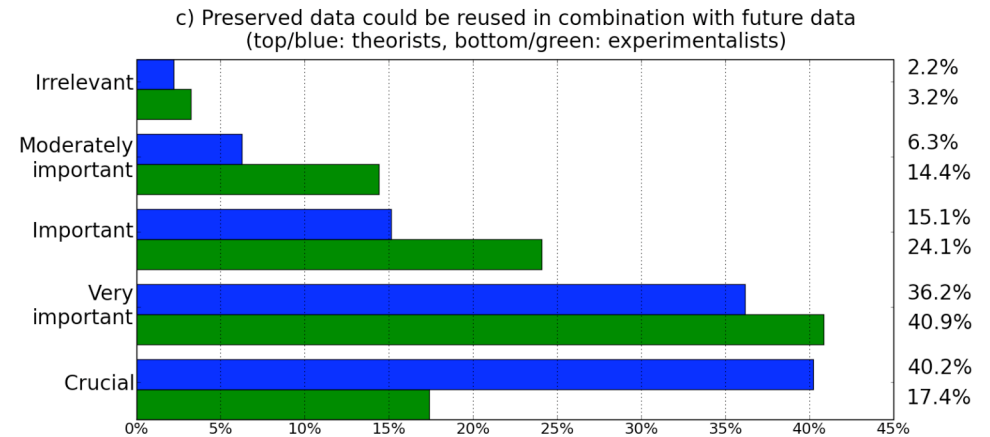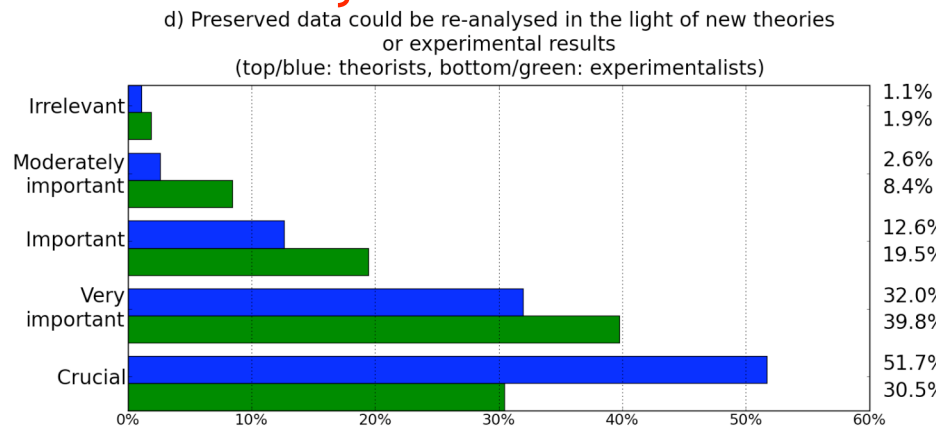(top/blue: theorists, bottom/green: experimentalists)



| | |
|---|---|
| Irrelevant | 0.4% |
| | 0.9% |
| Moderately important | 3.3% |
| | 8.7% |
| Important | 15.2% |
| | 25.6% |
| Very important | 41.7% |
| | 40.5% |
| Crucial | 39.5% |
| | 24.3% |

# The importance of preservation

## Future independent checks

b) Preserved data would allow future independent verification of results
(top/blue: theorists, bottom/green: experimentalists)

| | theorists | experimentalists |
|---|---|---|
| Irrelevant | 0.7% | 3.7% |
| Moderately important | 4.8% | 17.4% |
| Important | 13.7% | 24.8% |
| Very important | 37.3% | 32.0% |
| Crucial | 43.5% | 22.1% |

## Combine with future data

c) Preserved data could be reused in combination with future data
(top/blue: theorists, bottom/green: experimentalists)

| | theorists | experimentalists |
|---|---|---|
| Irrelevant | 2.2% | 3.2% |
| Moderately important | 6.3% | 14.4% |
| Important | 15.1% | 24.1% |
| Very important | 36.2% | 40.9% |
| Crucial | 40.2% | 17.4% |

## Re-analyse for future theories

d) Preserved data could be re-analysed in the light of new theories
or experimental results
(top/blue: theorists, bottom/green: experimentalists)

| | theorists | experimentalists |
|---|---|---|
| Irrelevant | 1.1% | 1.9% |
| Moderately important | 2.6% | 8.4% |
| Important | 12.6% | 19.5% |
| Very important | 32.0% | 39.8% |
| Crucial | 51.7% | 30.5% |

## Teaching and outreach

e) Preserved data could be used for teaching and outreach
(top/blue: theorists, bottom/green: experimentalists)

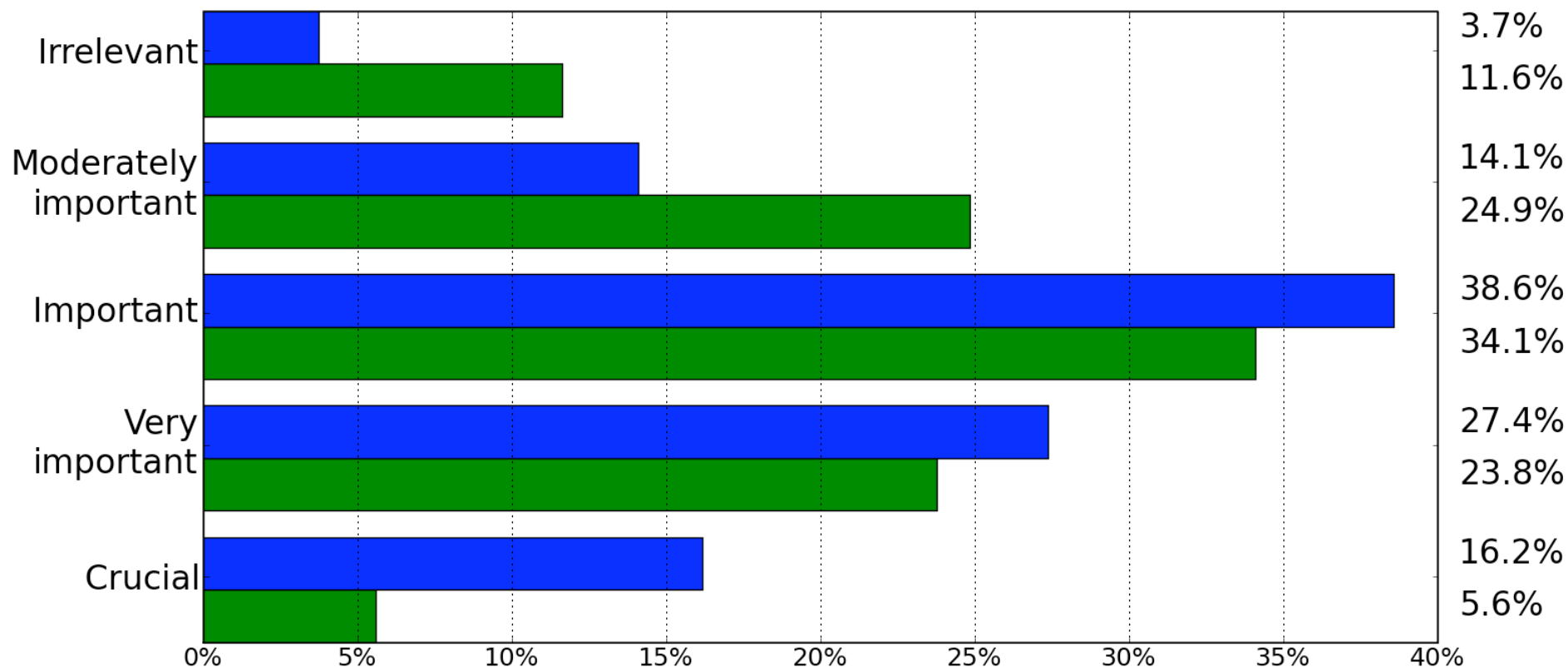| | theorists | experimentalists |
|---|---|---|
| Irrelevant | 9.9% | 8.3% |
| Moderately important | 27.9% | 32.6% |
| Important | 29.8% | 33.6% |
| Very important | 18.3% | 19.7% |
| Crucial | 14.1% | 5.8% |

# Why to preserve? - Compiling results

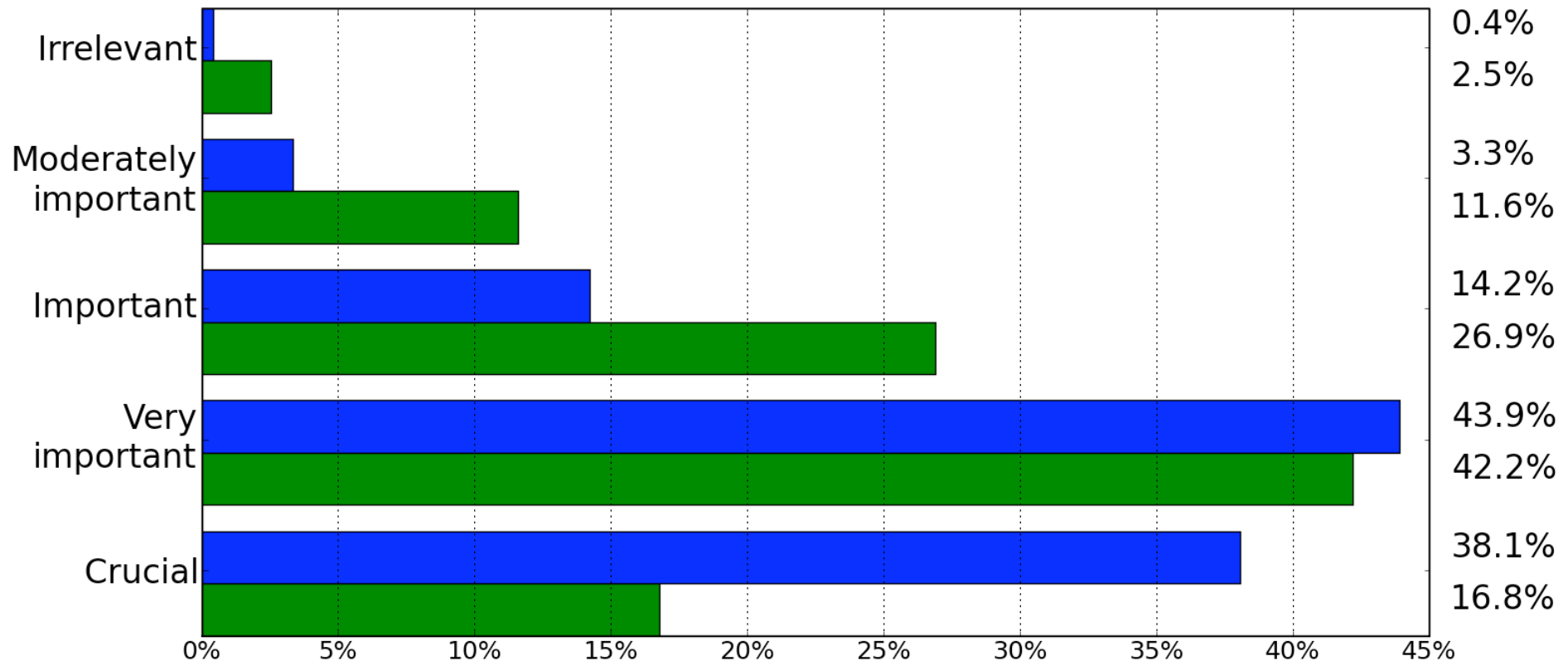How much importance would you attach to the following uses of preserved data ?

a) Compiling published results on a given subject (e.g. for a review)
(top/blue: theorists, bottom/green: experimentalists)



| | Theorists (blue) | Experimentalists (green) |
|---|---|---|
| Irrelevant | 3.7% | 11.6% |
| Moderately important | 14.1% | 24.9% |
| Important | 38.6% | 34.1% |
| Very important | 27.4% | 23.8% |
| Crucial | 16.2% | 5.6% |

# Why to preserve? - Testing new models

How much importance would you attach to the following uses of preserved data ?

b) Testing new models using preserved data
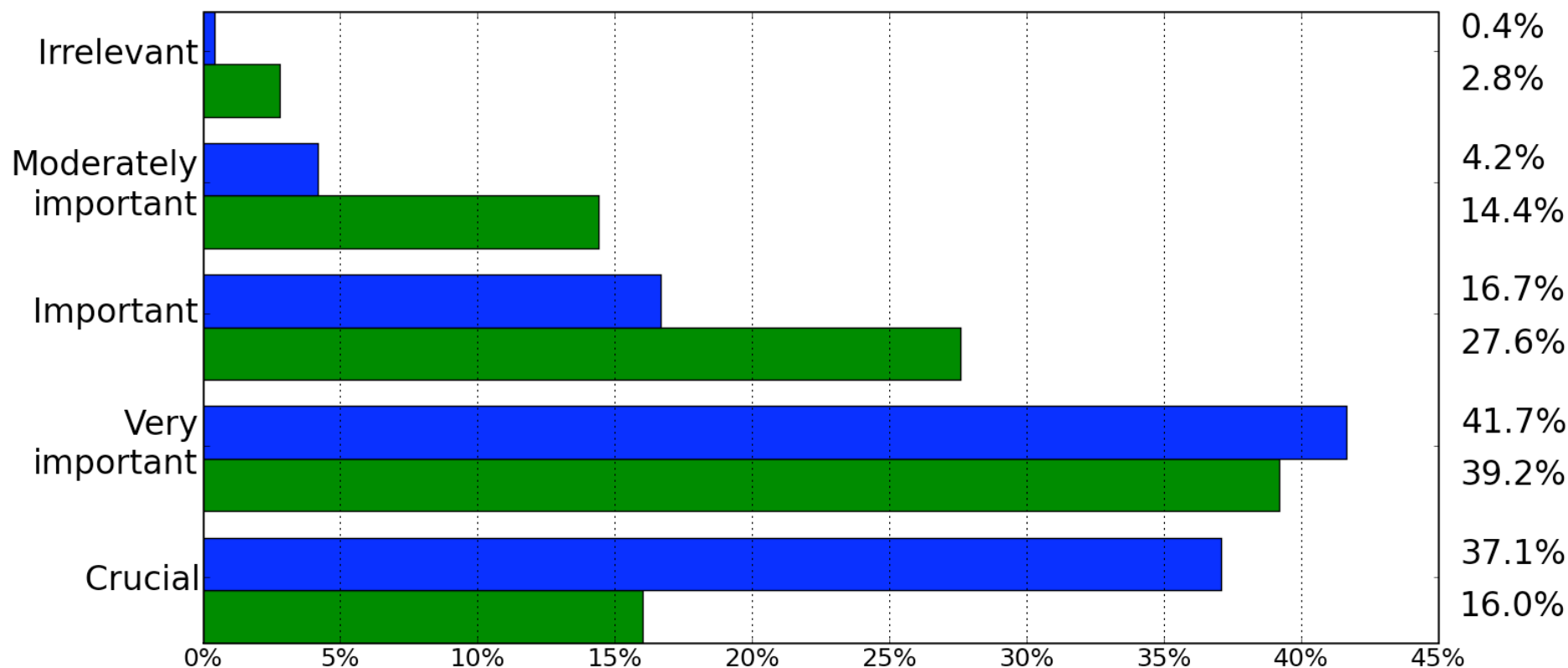(top/blue: theorists, bottom/green: experimentalists)



| | |
|---|---|
| Irrelevant | 0.4% |
| | 2.5% |
| Moderately important | 3.3% |
| | 11.6% |
| Important | 14.2% |
| | 26.9% |
| Very important | 43.9% |
| | 42.2% |
| Crucial | 38.1% |
| | 16.8% |

# Why to preserve? - Comparing past and future

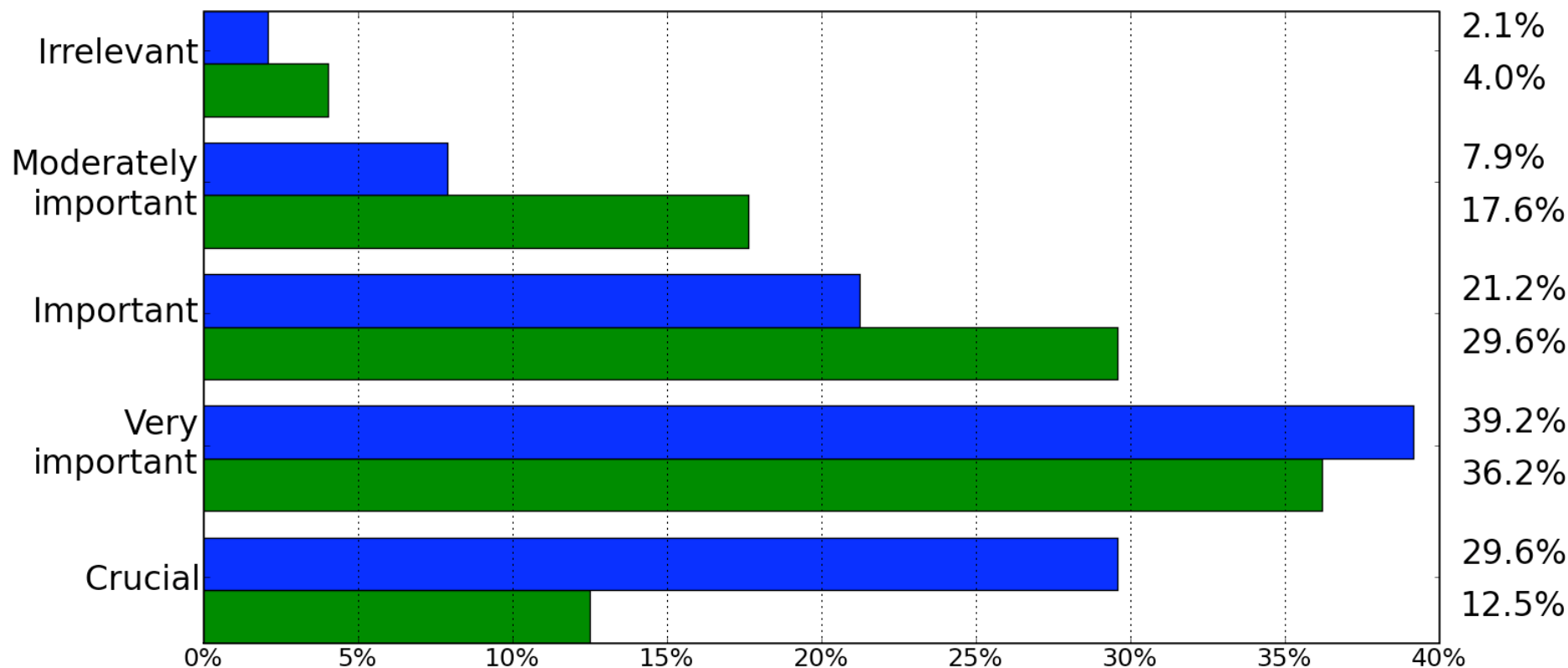How much importance would you attach to the following uses of preserved data ?

c) Showing compatibility of or detecting deviations between old and new experiments
(top/blue: theorists, bottom/green: experimentalists)



| | |
|---|---|
| Irrelevant | 0.4% / 2.8% |
| Moderately important | 4.2% / 14.4% |
| Important | 16.7% / 27.6% |
| Very important | 41.7% / 39.2% |
| Crucial | 37.1% / 16.0% |

# Why to preserve? - Combining past and future

How much importance would you attach to the following uses of preserved data ?

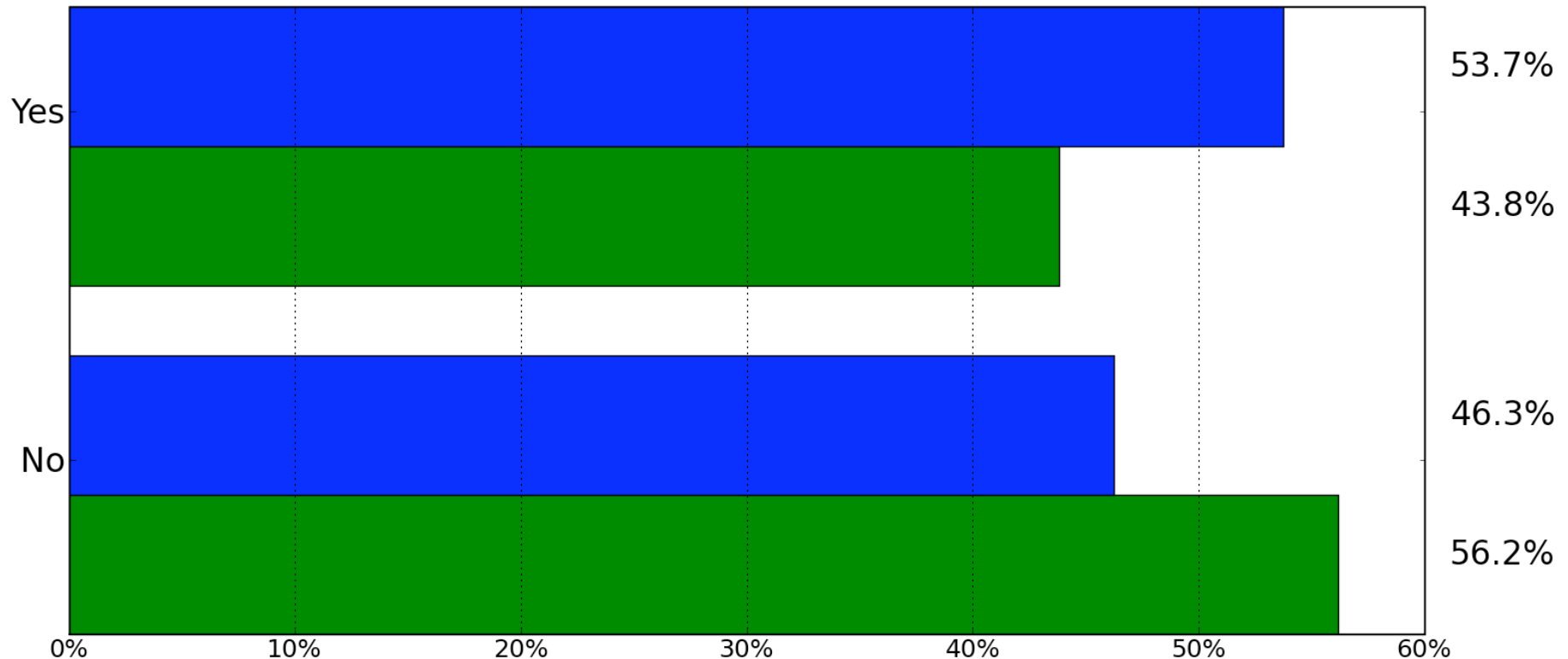d) Combining preserved data with new measurements
(top/blue: theorists, bottom/green: experimentalists)



| | Theorists (blue) | Experimentalists (green) |
|---|---|---|
| Irrelevant | 2.1% | 4.0% |
| Moderately important | 7.9% | 17.6% |
| Important | 21.2% | 29.6% |
| Very important | 39.2% | 36.2% |
| Crucial | 29.6% | 12.5% |

# Should we have started long ago?

Do you think that access to data from past experiments could
have improved your scientific results ?
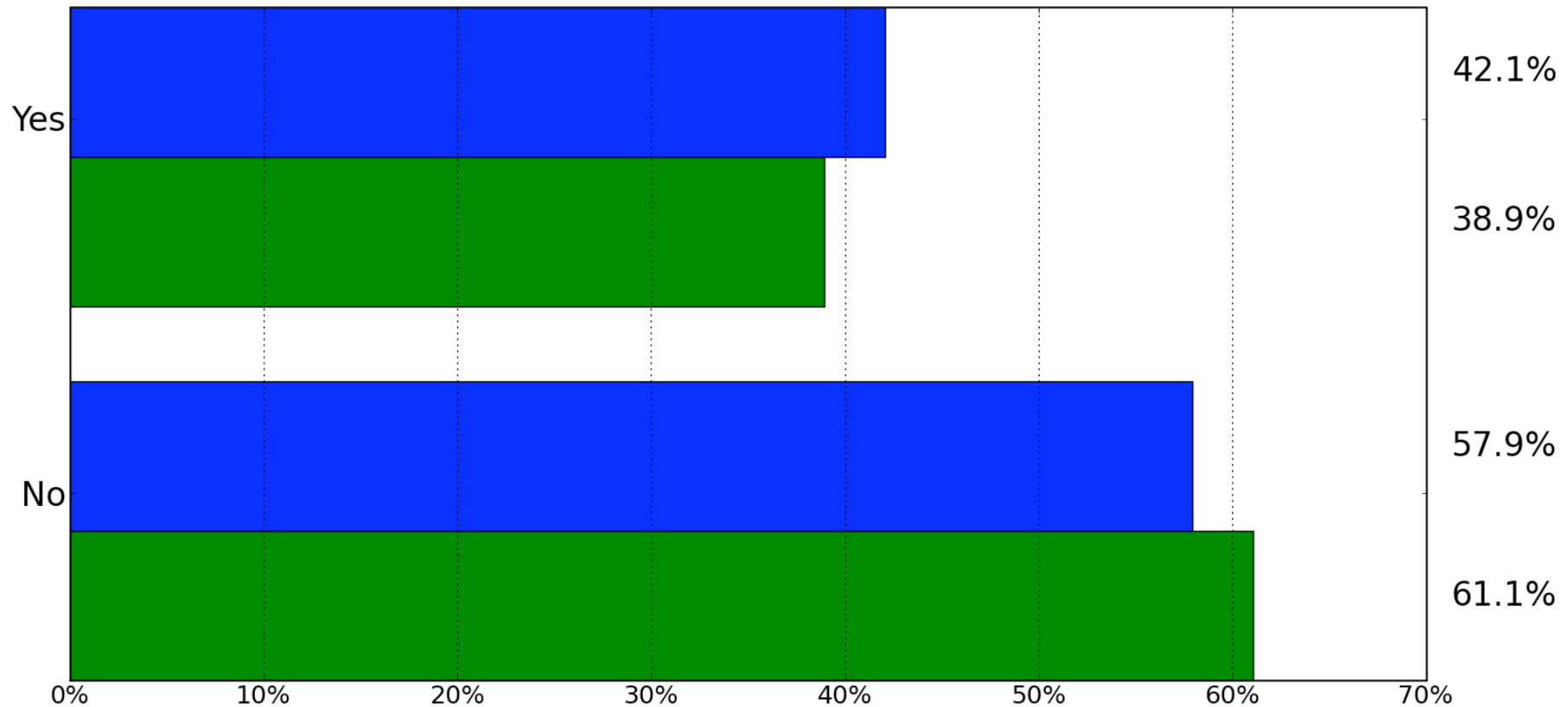(top/blue: theorists, bottom/green: experimentalists)



1. **Strong argument to plea for support to preserve**
2. **Demonstrate that preservation, re-use and (open) access cannot be divided**

# Did anything go wrong so far?



Do you think that in the past important HEP data have been lost ?
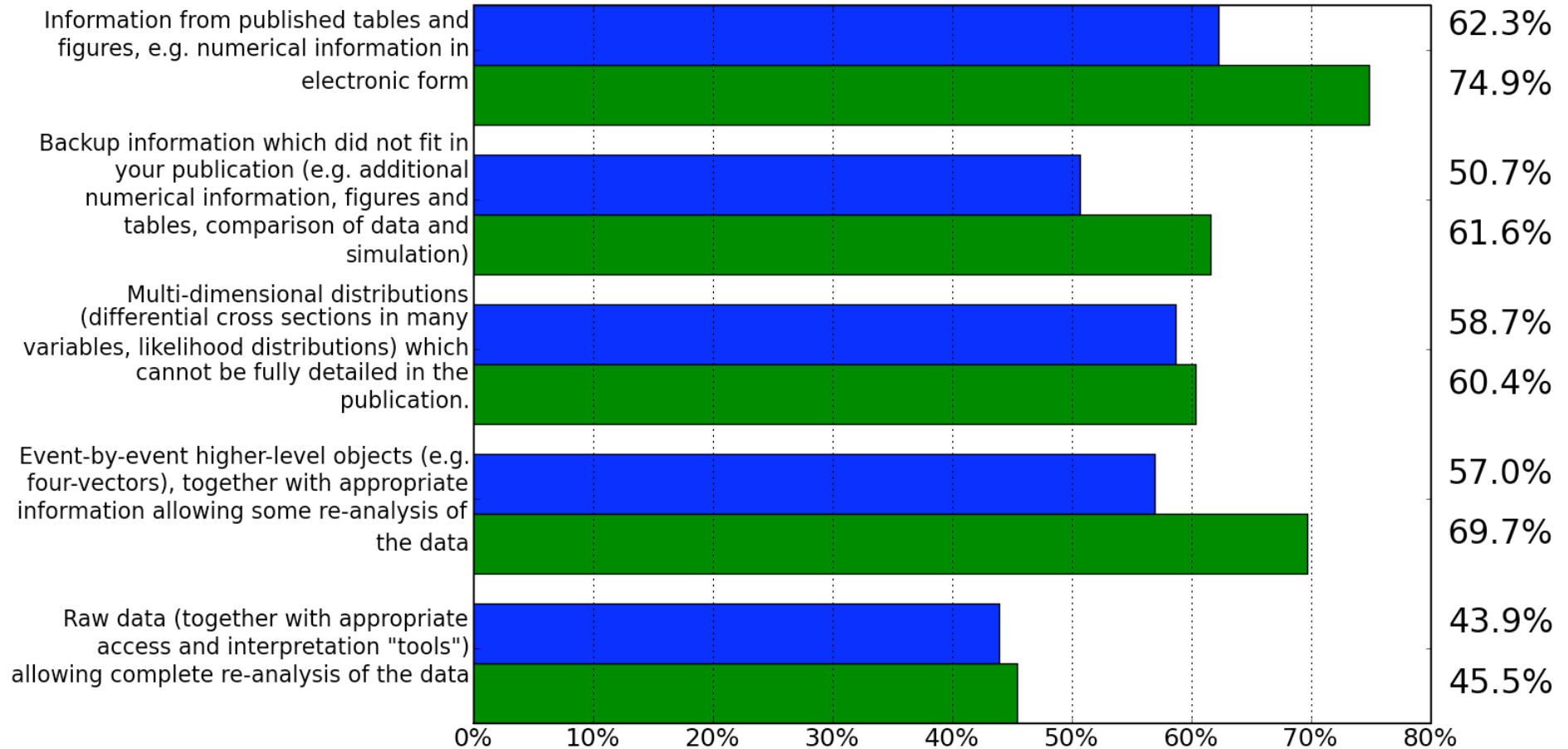(top/blue: theorists, bottom/green: experimentalists)

- Yes
  - 42.1%
  - 38.9%
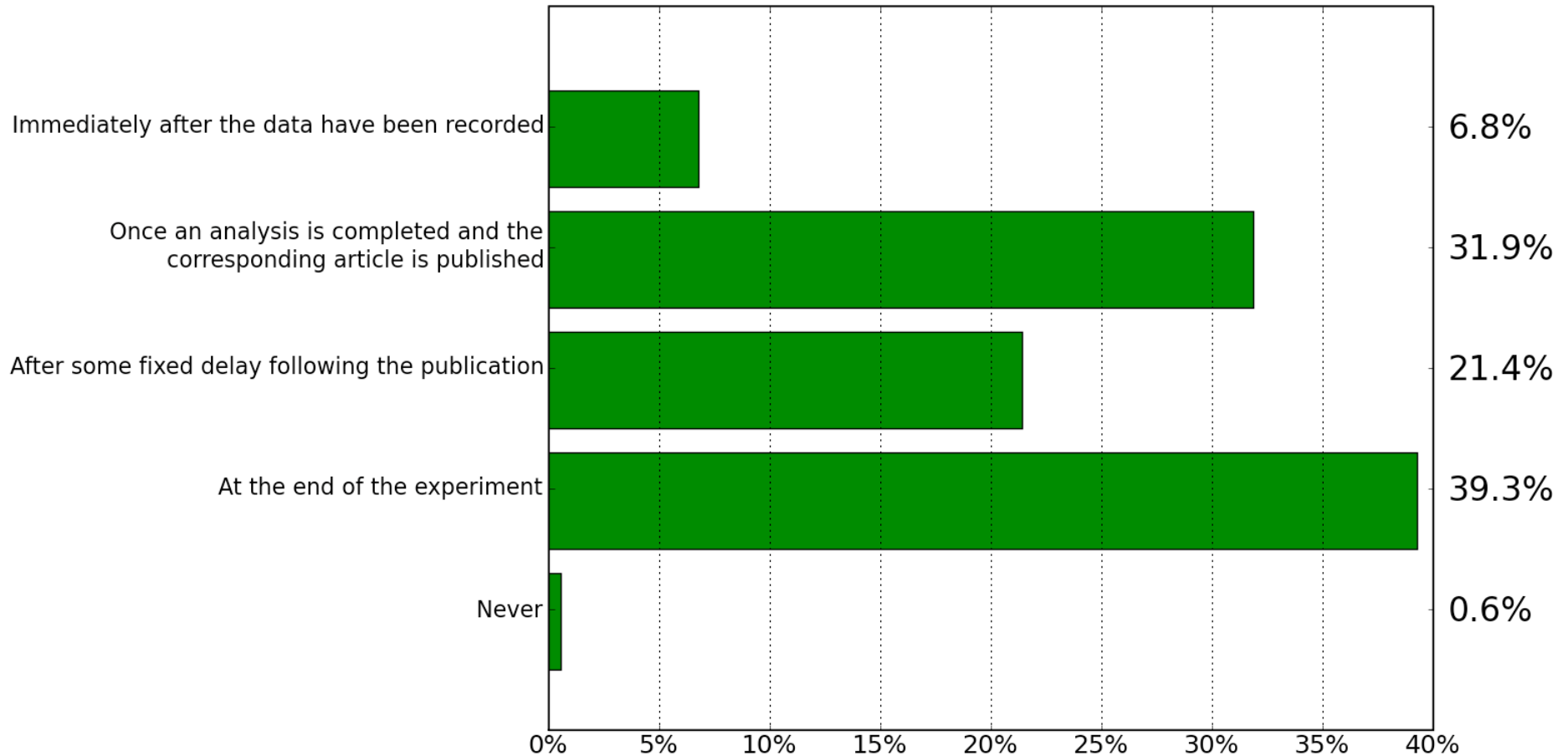- No
  - 57.9%
  - 61.1%

## Over optimistic? Over pessimistic?

# What to preserve?

At what level of detail should data be preserved ?
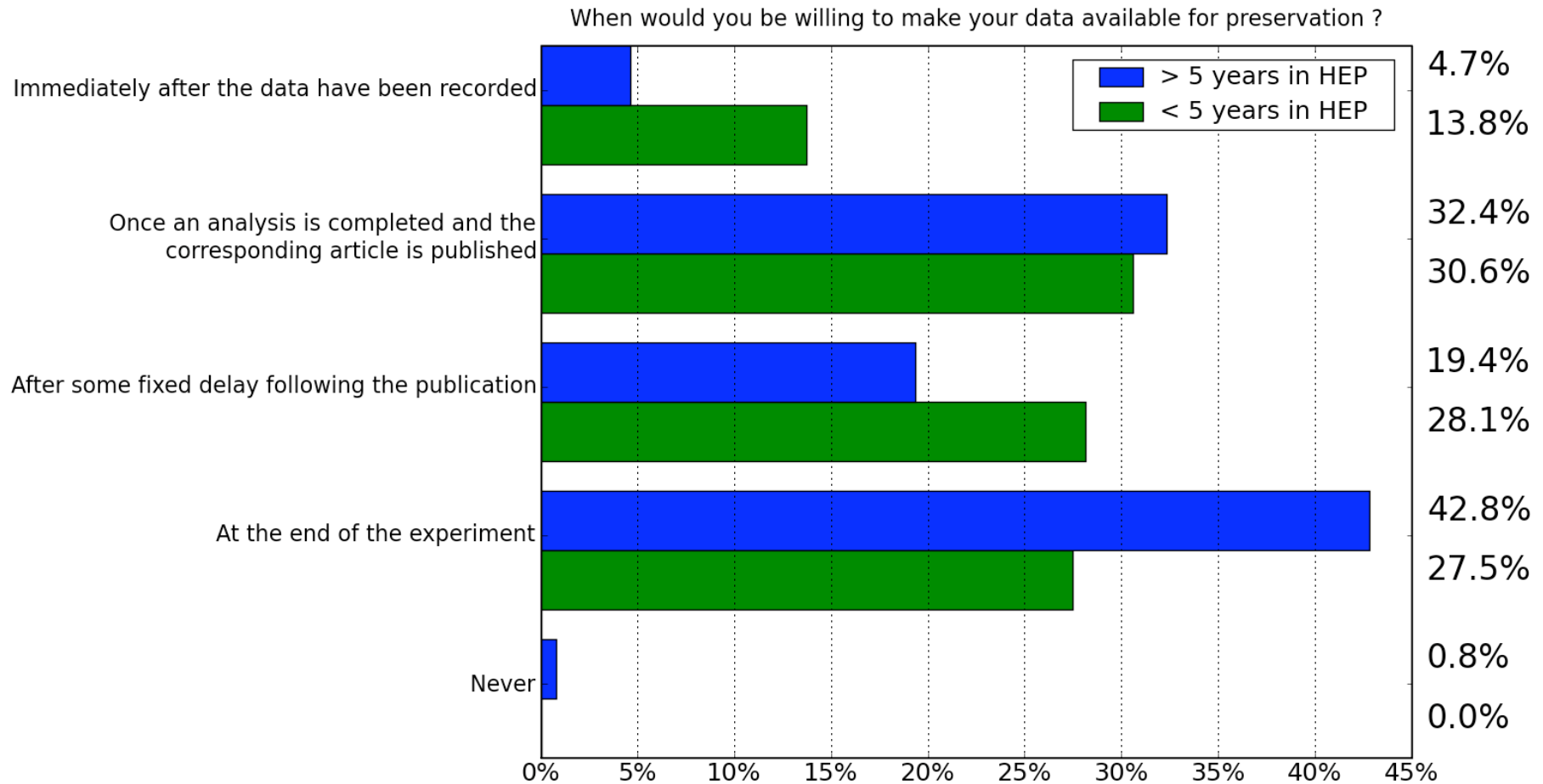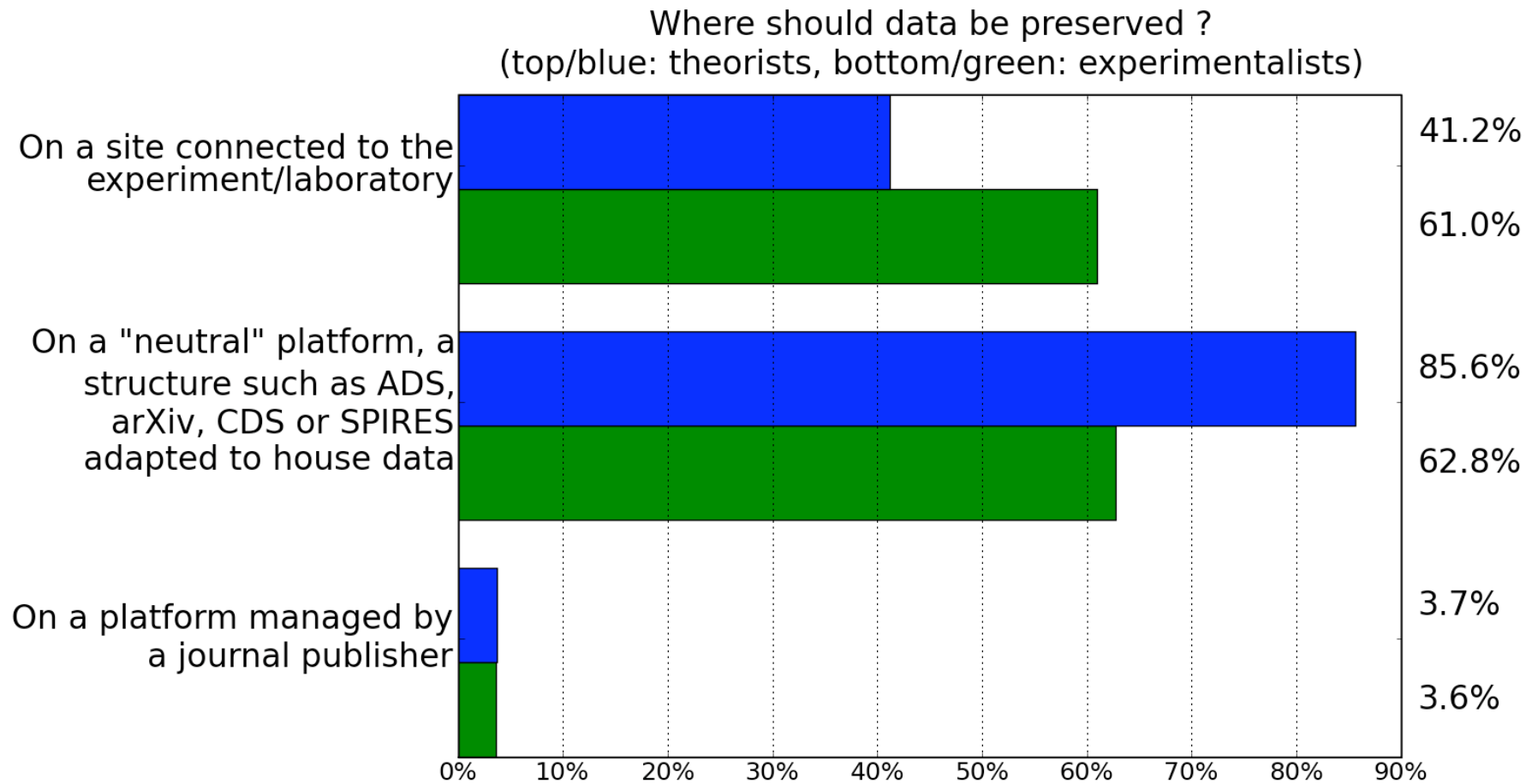(top/blue: theorists, bottom/green: experimentalists)

| Category | Theorists | Experimentalists |
|---|---|---|
| Information from published tables and figures, e.g. numerical information in electronic form | 62.3% | 74.9% |
| Backup information which did not fit in your publication (e.g. additional numerical information, figures and tables, comparison of data and simulation) | 50.7% | 61.6% |
| Multi-dimensional distributions (differential cross sections in many variables, likelihood distributions) which cannot be fully detailed in the publication. | 58.7% | 60.4% |
| Event-by-event higher-level objects (e.g. four-vectors), together with appropriate information allowing some re-analysis of the data | 57.0% | 69.7% |
| Raw data (together with appropriate access and interpretation "tools") allowing complete re-analysis of the data | 43.9% | 45.5% |

# When to preserve it?

When would you be willing to make your data available for preservation ?



| | |
|---|---|
| Immediately after the data have been recorded | 6.8% |
| Once an analysis is completed and the corresponding article is published | 31.9% |
| After some fixed delay following the publication | 21.4% |
| At the end of the experiment | 39.3% |
| Never | 0.6% |

# When to preserve it?



When would you be willing to make your data available for preservation ?

| | > 5 years in HEP | < 5 years in HEP |
|---|---|---|
| Immediately after the data have been recorded | 4.7% | 13.8% |
| Once an analysis is completed and the corresponding article is published | 32.4% | 30.6% |
| After some fixed delay following the publication | 19.4% | 28.1% |
| At the end of the experiment | 42.8% | 27.5% |
| Never | 0.8% | 0.0% |

# Where to preserve?



Where should data be preserved ?
(top/blue: theorists, bottom/green: experimentalists)

# Reality check #1: how though is it to preserve?

How much additional effort do you think is needed for the preservation of your data in a re-usable form
(in percent of the overall effort invested in the production and analysis of the data) ?
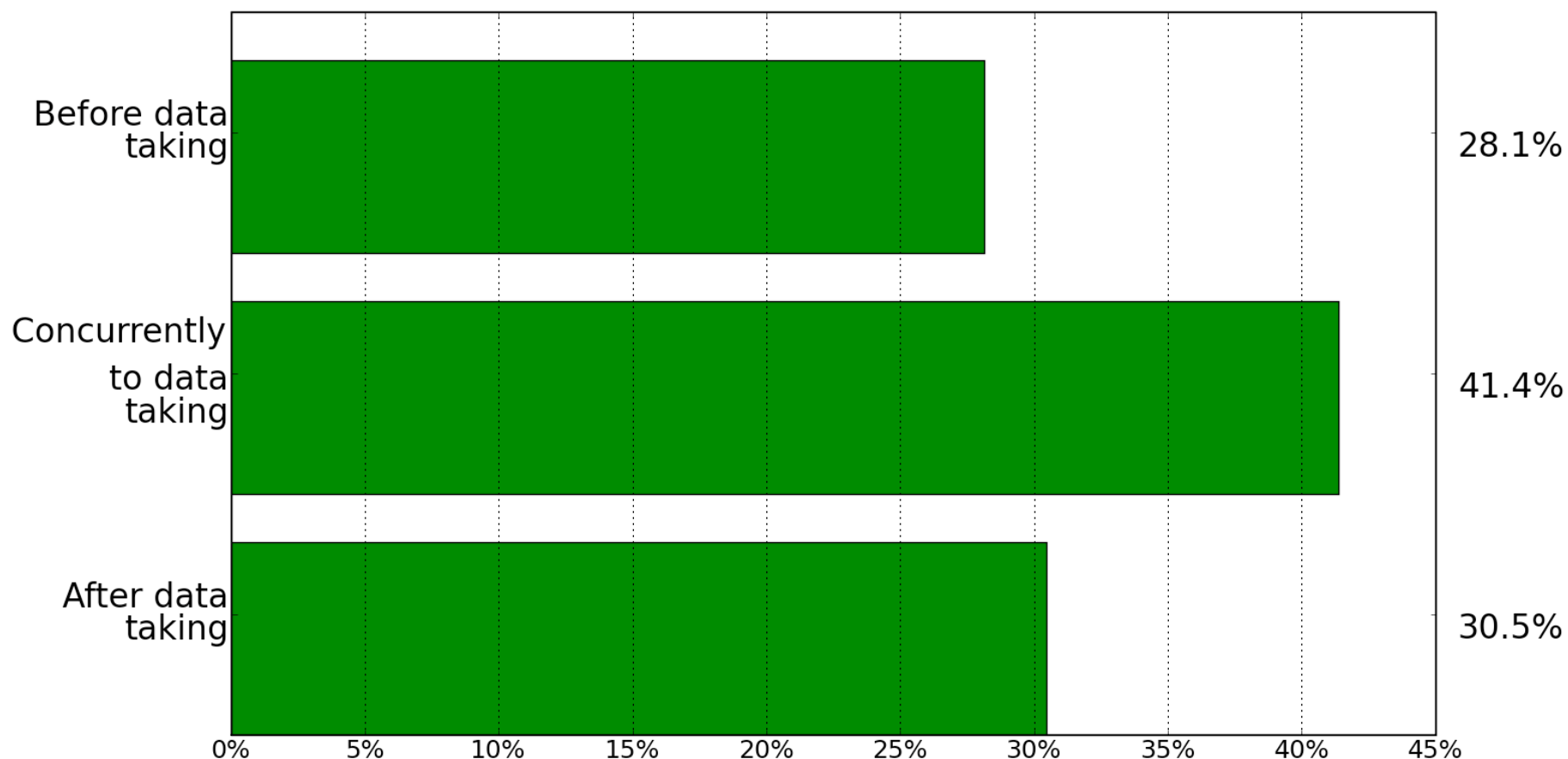
# Reality check #1: how though is it to preserve?

How much additional effort do you think is needed for the preservation of your data in a re-usable form
(in percent of the overall effort invested in the production and analysis of the data) ?
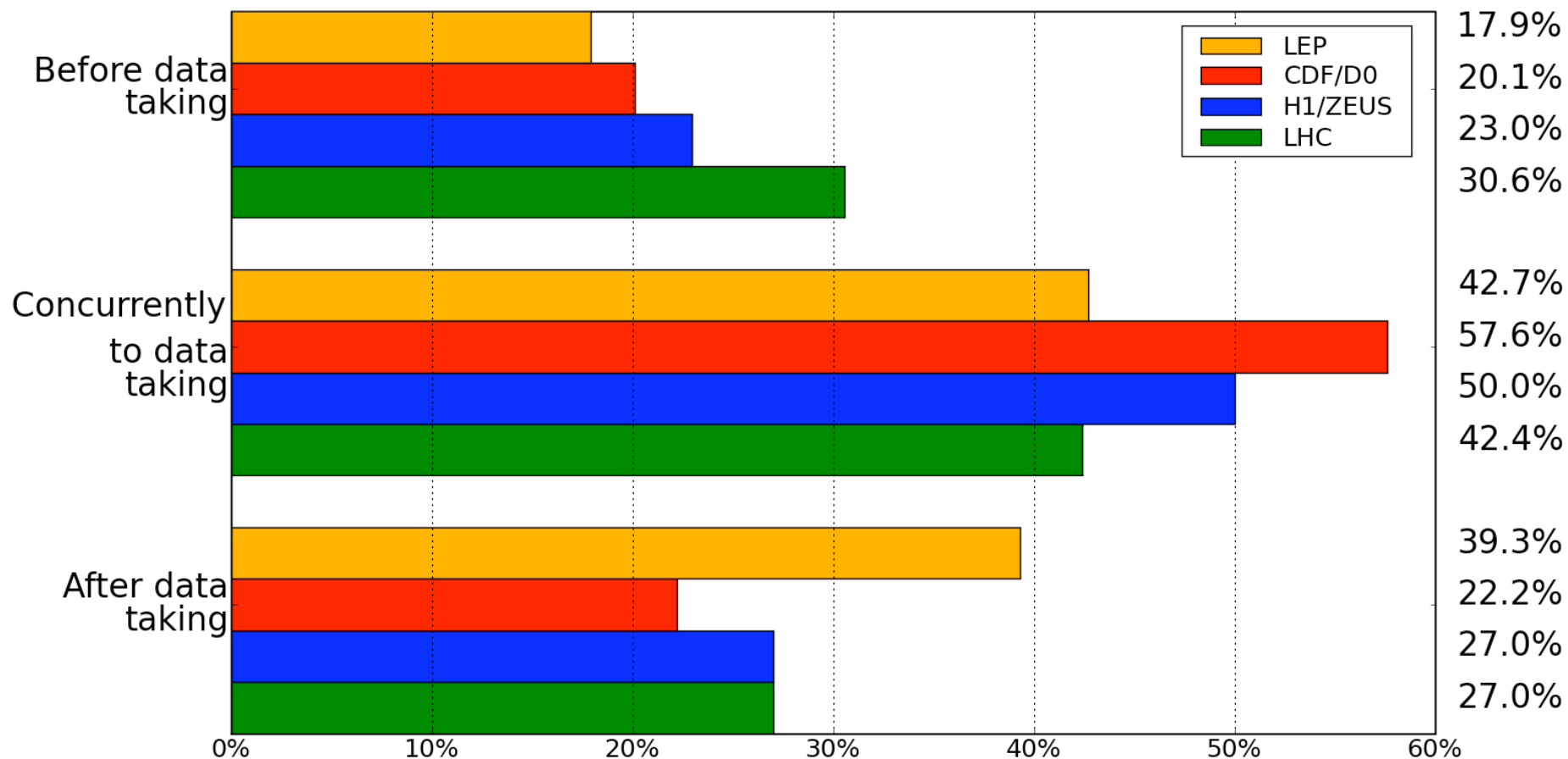


Legend:
- LEP
- CDF/D0
- H1/ZEUS
- LHC

< 1%
- 6.7%
- 6.8%
- 5.3%
- 6.0%

1-10 %
- 47.9%
- 50.7%
- 46.7%
- 50.7%

10-50%
- 37.0%
- 32.4%
- 42.7%
- 36.9%

> 50%
- 8.4%
- 10.1%
- 5.3%
- 6.4%

# Reality check #2: when to preserve?

In your opinion, when should this effort start in order to be the most effective ?



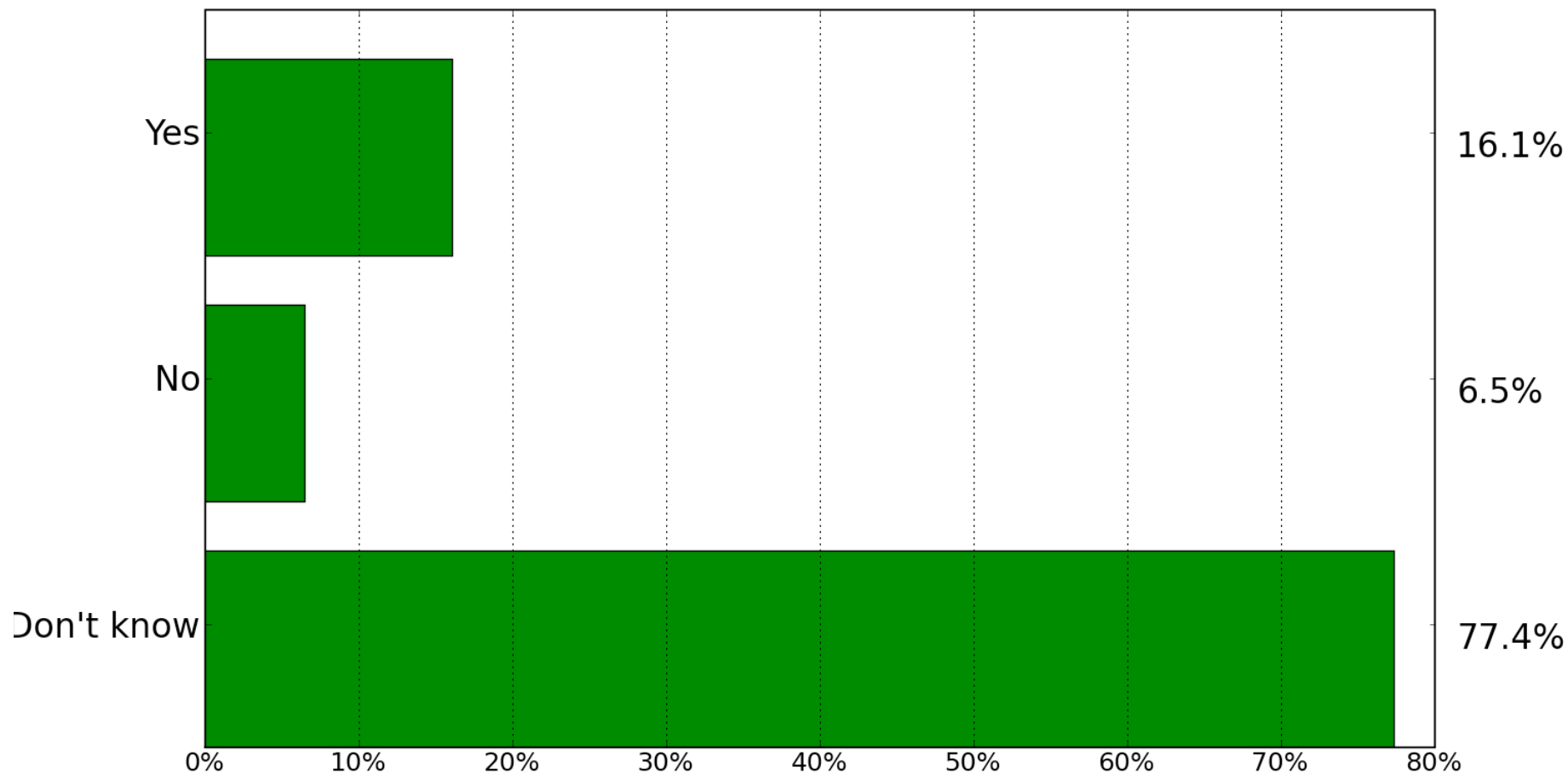| | |
|---|---|
| Before data taking | 28.1% |
| Concurrently to data taking | 41.4% |
| After data taking | 30.5% |

# Reality check #2: when to preserve?

In your opinion, when should this effort start in order to be the most effective ?

# Reality check #3: is it doable?
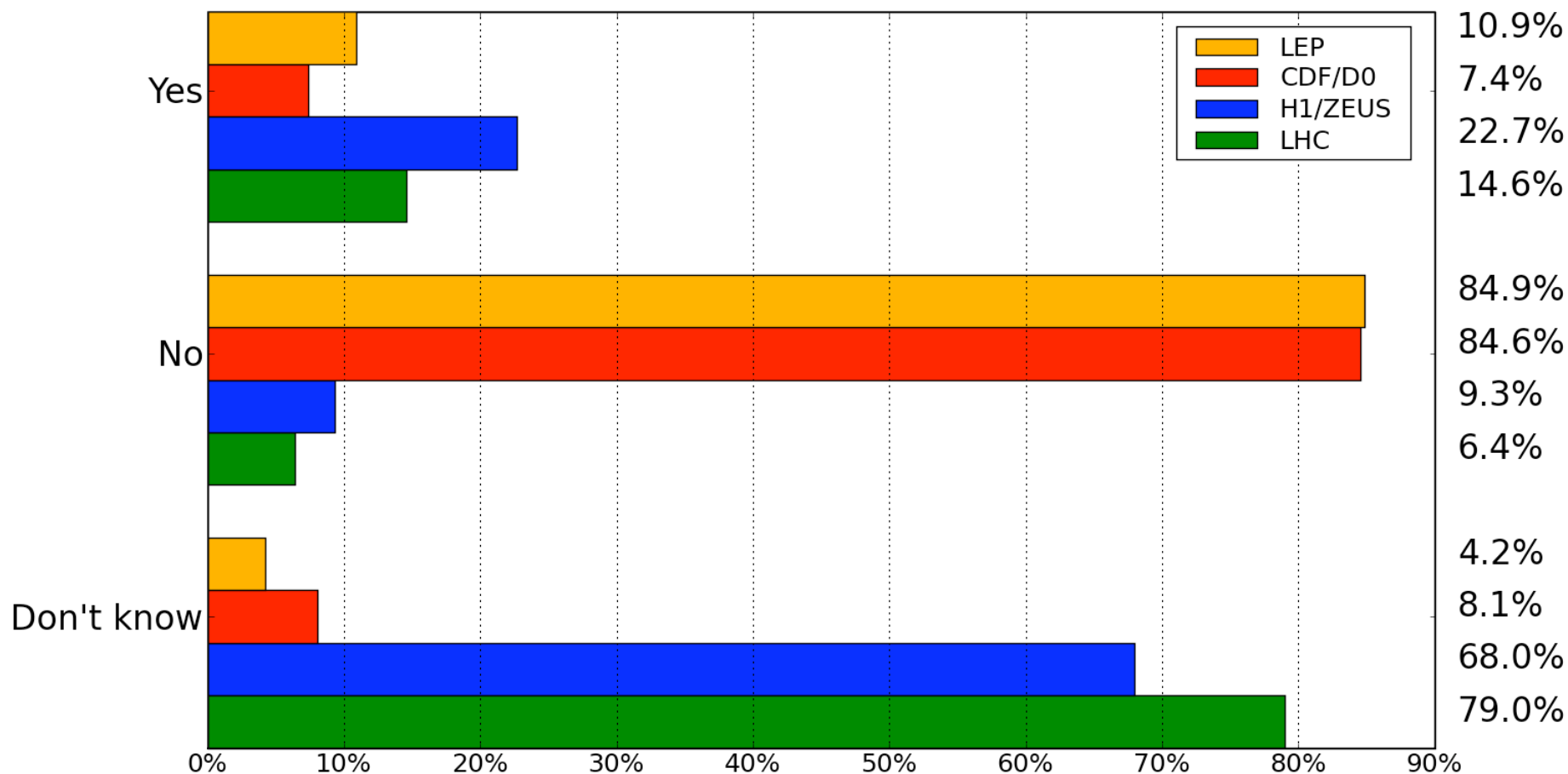


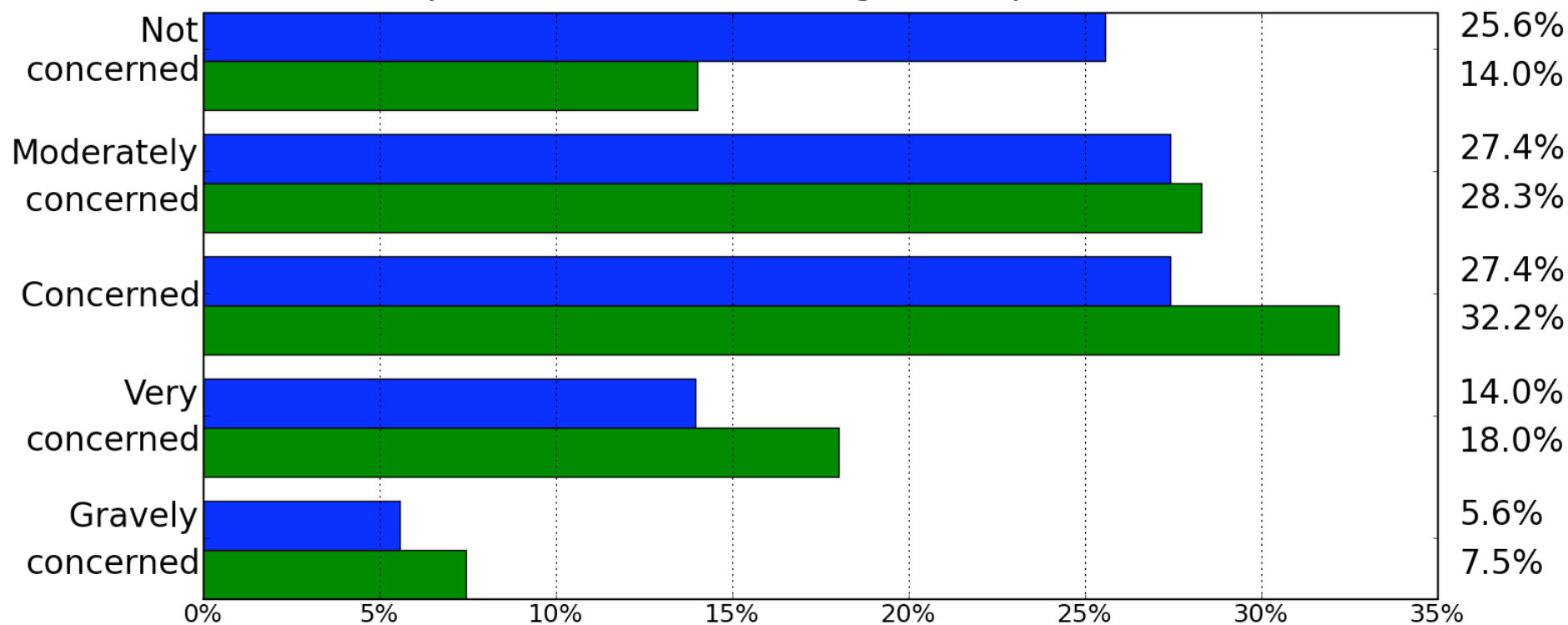Will your experiment/collaboration/organisation be able to invest this effort ?

- Yes — 16.1%
- No — 6.5%
- Don't know — 77.4%

# Reality check #3: is it doable?

Will your experiment/collaboration/organisation be able to invest this effort ?



| | |
|---|---|
| **Yes** | 10.9% (LEP) |
| | 7.4% (CDF/D0) |
| | 22.7% (H1/ZEUS) |
| | 14.6% (LHC) |
| **No** | 84.9% (LEP) |
| | 84.6% (CDF/D0) |
| | 9.3% (H1/ZEUS) |
| | 6.4% (LHC) |
| **Don't know** | 4.2% (LEP) |
| | 8.1% (CDF/D0) |
| | 68.0% (H1/ZEUS) |
| | 79.0% (LHC) |

Legend: LEP, CDF/D0, H1/ZEUS, LHC

# Ideal-case worries: getting credit

To what extent are you concerned about the following issues
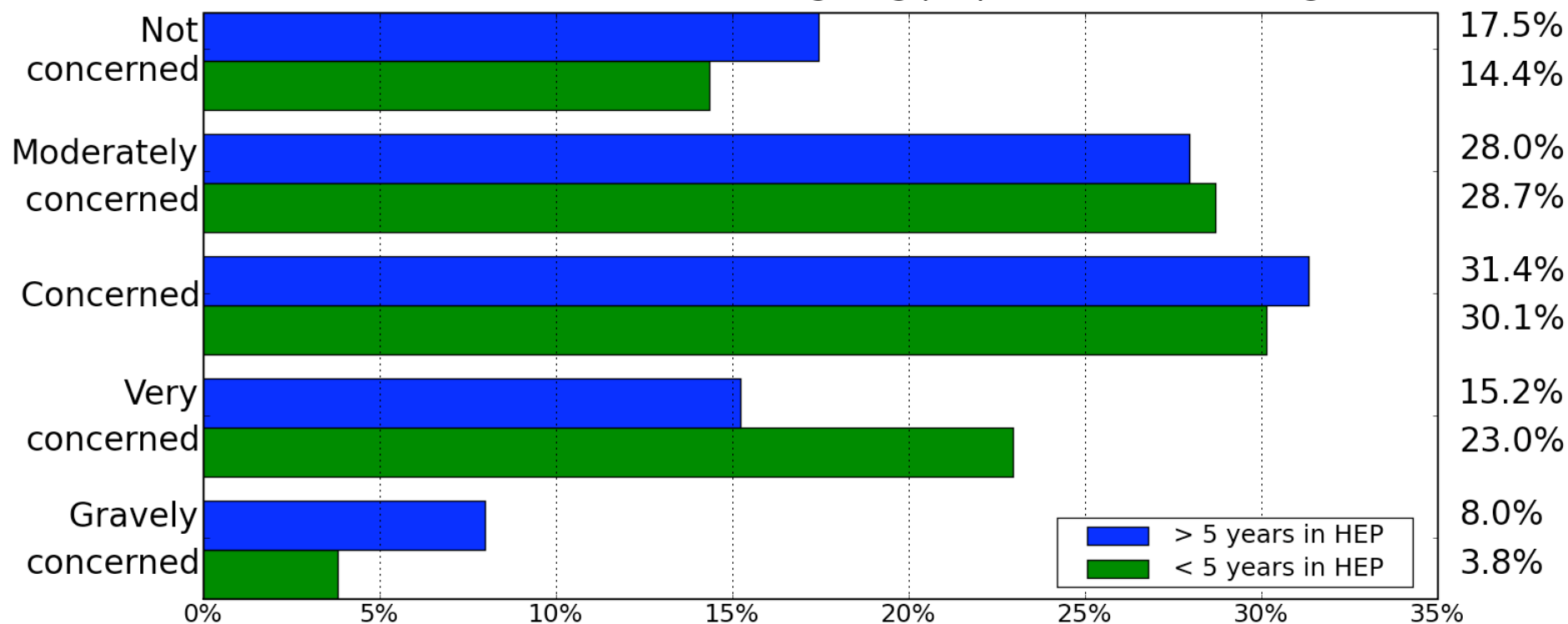related to giving access to preserved data ?

a) Preserved data could be used without giving proper credit to the original authors
(top/blue: theorists, bottom/green: experimentalists)

# Ideal-case worries: getting credit

To what extent are you concerned about the following issues
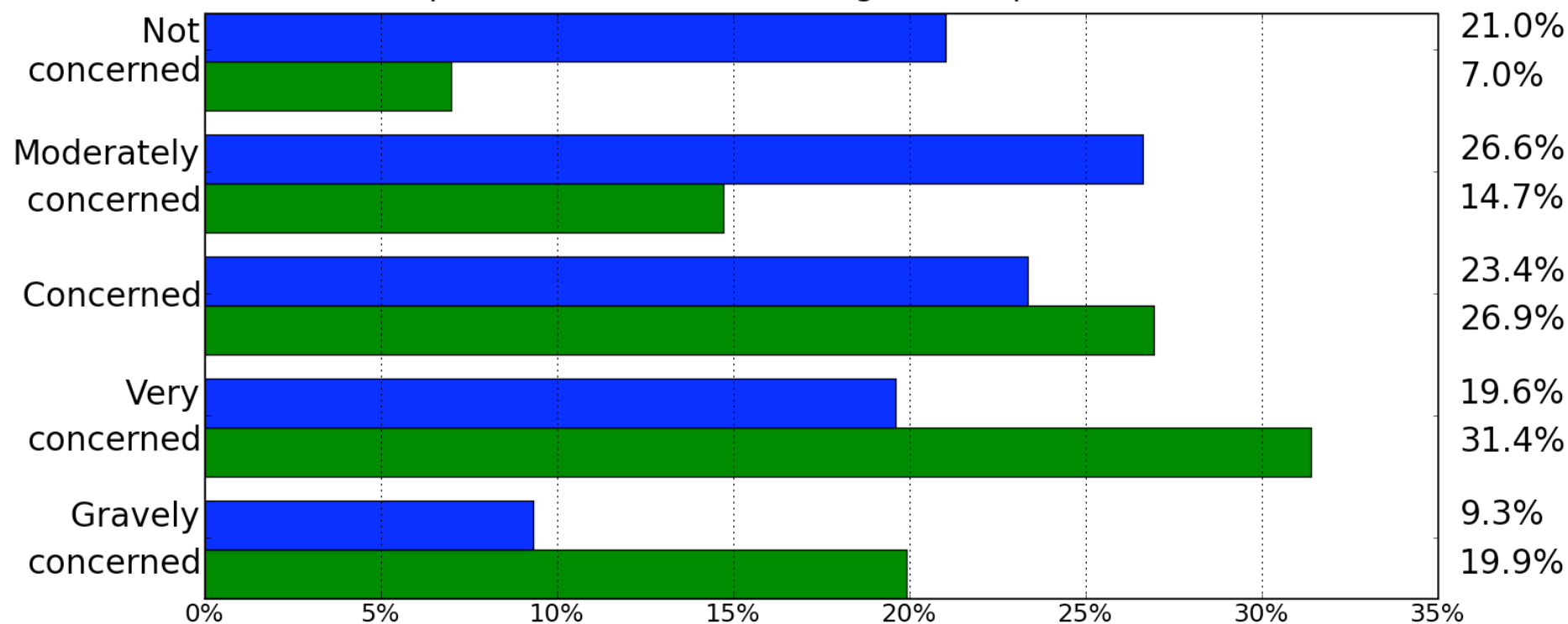related to giving access to preserved data ?

a) Preserved data could be used without giving proper credit to the original authors



| | > 5 years in HEP | < 5 years in HEP |
|---|---|---|
| Not concerned | 17.5% | 14.4% |
| Moderately concerned | 28.0% | 28.7% |
| Concerned | 31.4% | 30.1% |
| Very concerned | 15.2% | 23.0% |
| Gravely concerned | 8.0% | 3.8% |

# Ideal-case worries: inflation/noise

To what extent are you concerned about the following issues
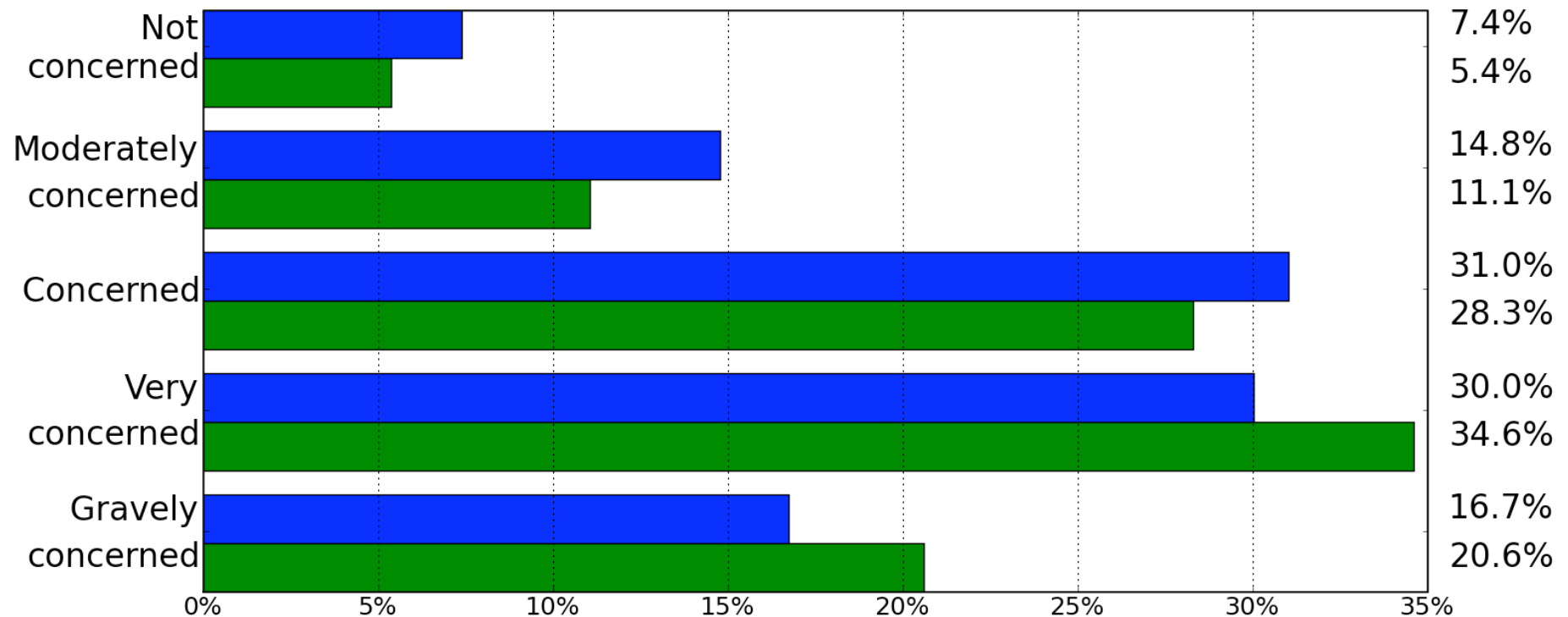related to giving access to preserved data ?

b) Uncontrolled access to data may lead to an inflation of incorrect results
(top/blue: theorists, bottom/green: experimentalists)



| Category | Theorists (blue) | Experimentalists (green) |
|---|---|---|
| Not concerned | 21.0% | 7.0% |
| Moderately concerned | 26.6% | 14.7% |
| Concerned | 23.4% | 26.9% |
| Very concerned | 19.6% | 31.4% |
| Gravely concerned | 9.3% | 19.9% |

# Ideal-case worries: documentation

If you were to re-use preserved data, to what extent would
you be concerned by the following scenarios ?

d) I am not using the data correctly
(top/blue: theorists, bottom/green: experimentalists)



| | |
|---|---|
| Not concerned | 7.4% / 5.4% |
| Moderately concerned | 14.8% / 11.1% |
| Concerned | 31.0% / 28.3% |
| Very concerned | 30.0% / 34.6% |
| Gravely concerned | 16.7% / 20.6% |

# PARSE.Insight HEP Case Study: next steps

- Many free-text open-ended questions still to analyse concerning: threats, opportunities, feasibility, access regulation, trustworthiness... Compile relevant excerpts and quantification of results

- 182 respondent (15%) made themselves available for an interview to express more opinion. Identify relevant subjects and run interviews

- Go to known "superusers" with same interview questions

## Any inputs?