

“The HERA Inheritance”.

The DP project at DESY, associated costs and lessons learned

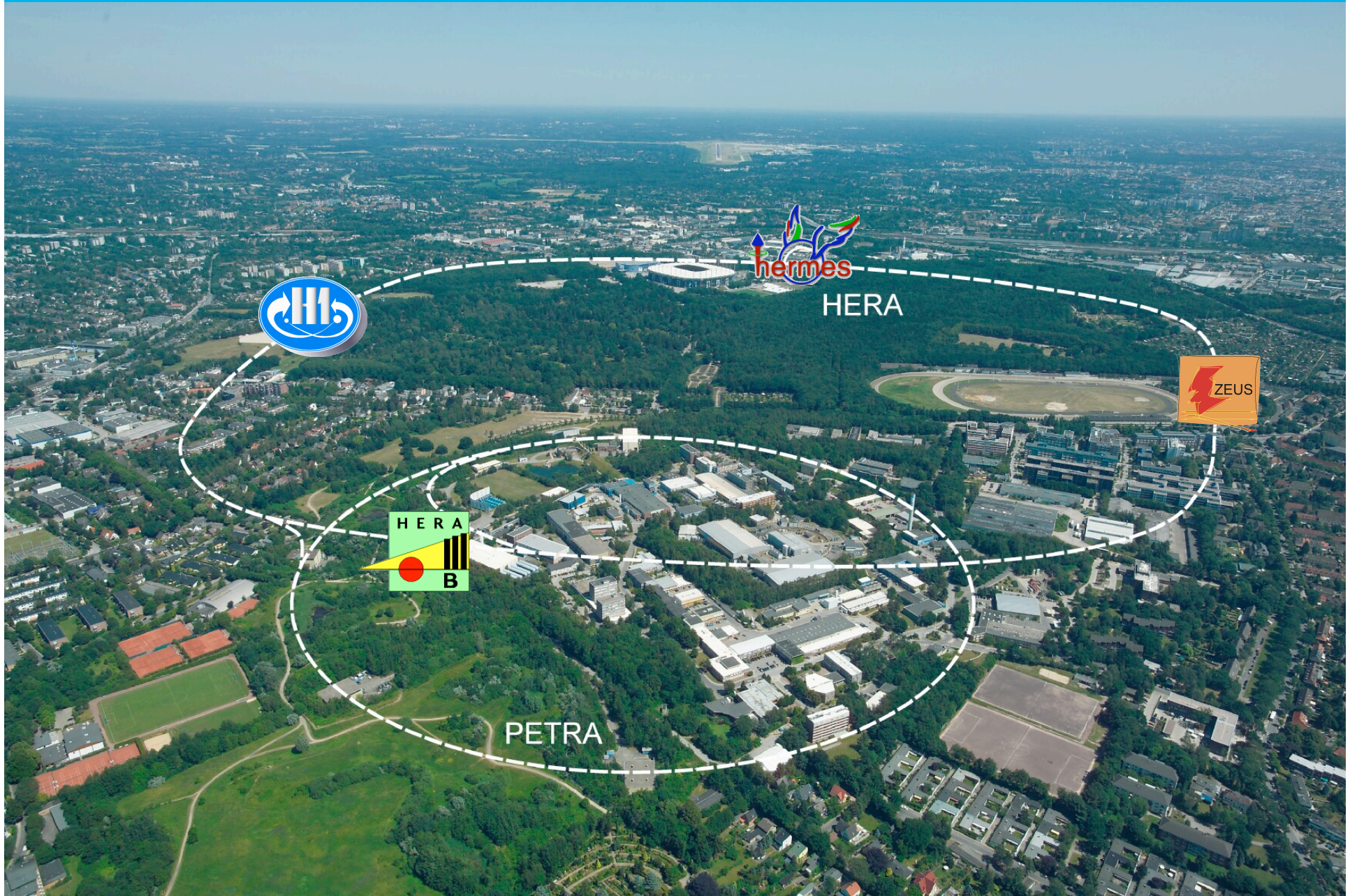


David South (DESY)

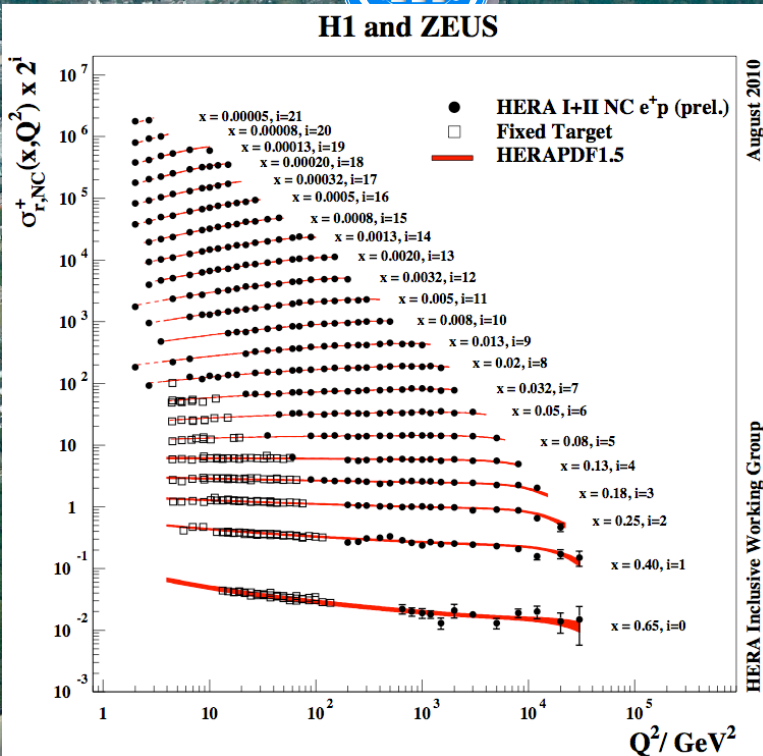
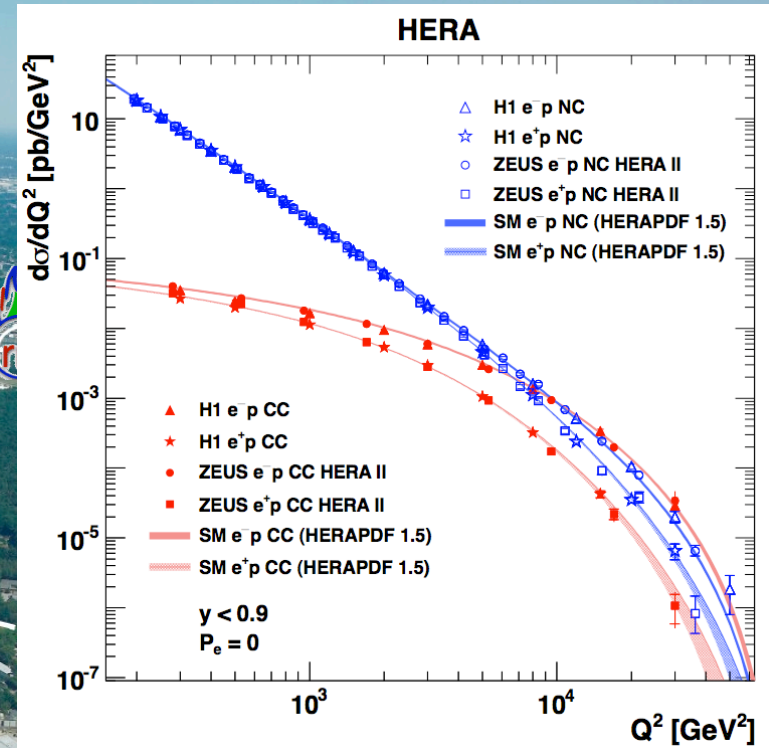
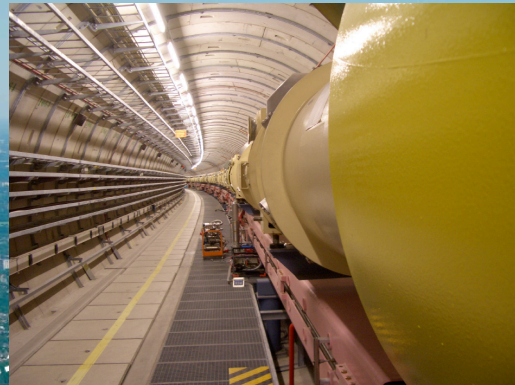
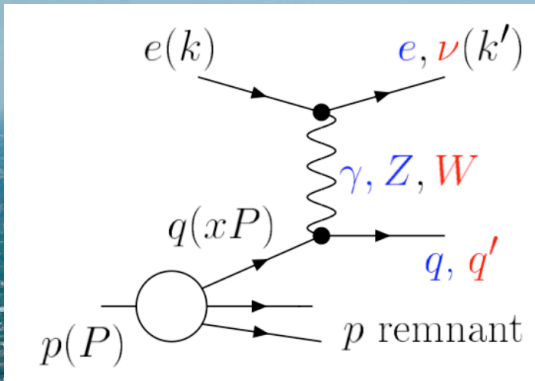
8th DPHEP Workshop: “Costs of Curation”
CERN, 13th of January 2014



1992: Hadron-Electron Ring Accelerator (HERA) @ DESY



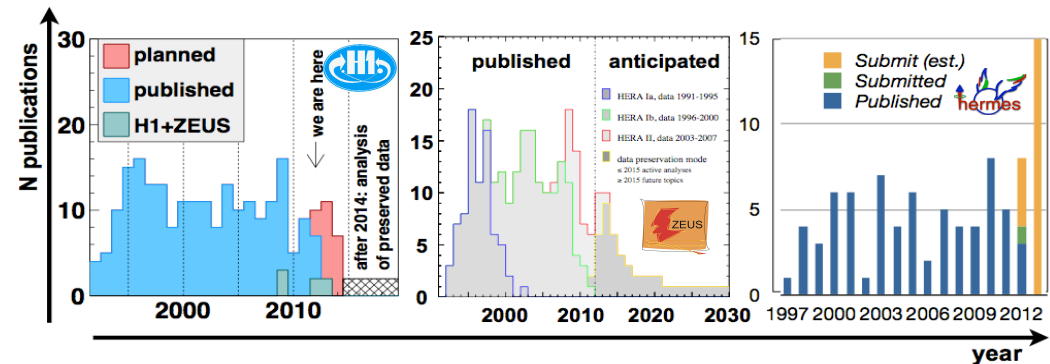
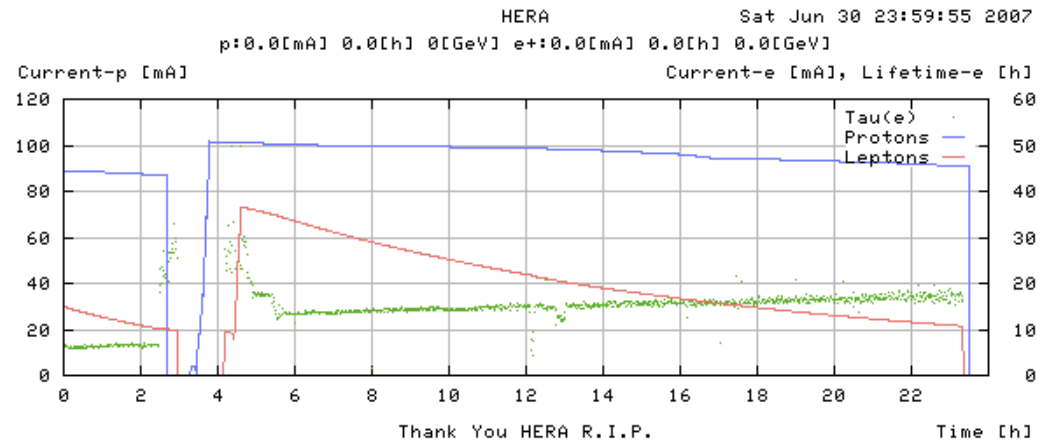
1992: Hadron-Electron Ring Accelerator (HERA) @ DESY



- The world's only electron-proton collider, collisions at H1 and ZEUS 1994-2007
- Precise picture of the proton, crucial measurements for hadron colliders
- Many other unique physics measurements

End of data taking at HERA: June 30th 2007

- > Unique period of time in HEP history: change from many running experiments of various types to essentially only one
- > HERA, stopped taking data 6.5 years ago – so what's happened since then?
- > Much like LEP before us and seen by BaBar, publications still continue well after data taking: *~25% of total so far!*
 - H1: 55 papers since June 2007, out of a total of 218
 - ZEUS: 64 out of a total of 241



DPHEP activity at DESY since 2008

- The first few years after data taking: 2008-2010
 - Formation of initial ideas, first DPHEP workshops
 - Grand surveys done: data, hardware, software, technologies
 - Establishing the physics case for data preservation
 - Defining the DPHEP preservation levels: HERA experiments plan for level 3-4
 - Finding the people to do the work



DPHEP activity at DESY since 2008

> The first few years after data taking: 2008-2010

- Formation of initial ideas, first DPHEP workshops
- Grand surveys done: data, hardware, software, technologies
- Establishing the physics case for data preservation
- Defining the DPHEP preservation levels: HERA experiments plan for level 3-4
- Finding the people to do the work

> Key areas of activity at DESY since 2011

1. Preparation of the data for preservation and archival storage of the data themselves
2. Data preservation: really preservation of software + environment: the `sp-system`
3. Documentation: INSPIRE, digital meta-data and non-digital material
4. Governance, future collaboration structures and open access/public data, outreach



DPHEP activity at DESY since 2008

> The first few years after data taking: 2008-2010

- Formation of initial ideas, first DPHEP workshops
- Grand surveys done: data, hardware, software, technologies
- Establishing the physics case for data preservation
- Defining the DPHEP preservation levels: HERA experiments plan for level 3-4
- Finding the people to do the work

> Key areas of activity at DESY since 2011

1. Preparation of the data for preservation and archival storage of the data themselves
2. Data preservation: really preservation of software + environment: the `sp-system`
3. Documentation: INSPIRE, digital meta-data and non-digital material
4. Governance, future collaboration structures and open access/public data, outreach

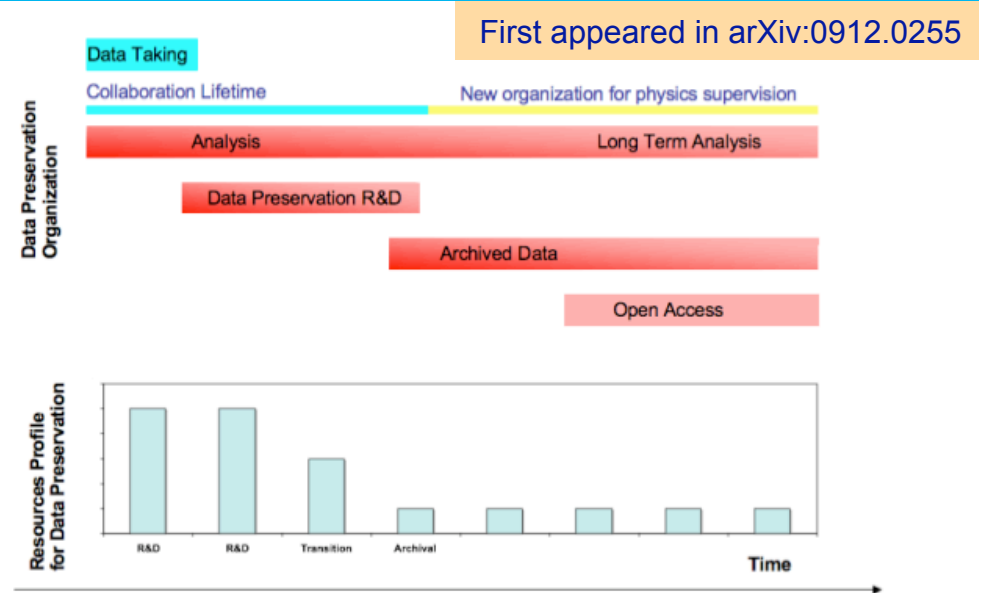
> Some things have been completed, typically well defined tasks such as the documentation (rather specialised person-power), other things still on going, including final preparation of archival data storage



The DESY-DPHEP Group

- > During first years, regularly more than a dozen people involved
- > Group made up of people from H1, ZEUS and HERMES as well as DESY-IT and DESY-Library
 - The available person-power has declined, in line with the model presented in the first DPHEP publication
 - 2014: Now only a couple of people involved

		2011	2012	2013	Translates into Position	2014 ⁺⁺
DESY-IT	Validation	1.0		0.5	3 year FTE 2011 – 2013	(0.5)
	Storage		1.0	0.5		
H1	Validation	0.5	1.0	0.5	2 year extension for 2011 – 2013	(0.5)
	Documentation	0.5	0.5		1 year extension for 2012	
ZEUS	Validation	0.5	1.0	0.5	(Initial) 2 year FTE 2011 – 2013	(0.5)
	Documentation	0.5	0.5		1 year FTE 2011 – 2012	
HERMES	Validation			0.5	0.5 year extension for 2013	(0.5)
	Documentation	0.5	0.5		1 year FTE 2011 – 2012	



- > Initial person-power estimates included provision for support in 2014 and beyond

- Long term support has proven difficult to secure, especially when trying to find the right people for the job
- All current DP person-power runs out this year. IT part will be covered for a further 2 years



Key area 1: Data for preservation and archival storage

- > Deciding which data (and MC) are needed for the long term depends on the preservation model assumed: level 4 goes back to the raw data
- > Final production of HERA data for preservation only completed last year; majority of MC production expected to be concluded this year

- > Estimates for final **DPHEP dataset** volume ready (including MC samples)

- Plan calls for two tape copies and an “always online” (disk) component
- Data which should be archived, but not online all the time: re-pack into larger files
- Costs not prohibitive on data volume basis

Expt	Online (TB)	Total (TB)
H1	250	500
ZEUS	250	250
HERMES	100	300
Total	600	1050
HERA-B	?	300

Different strategies visible

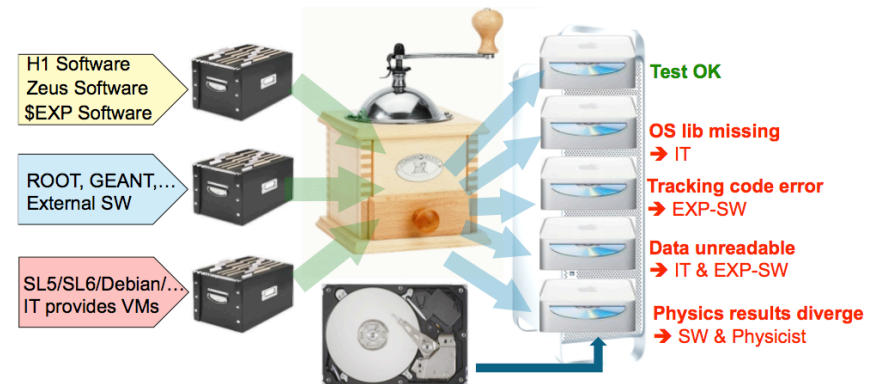
- > Dedicated system too costly in both hardware and support required
 - All collaborations use dCache for mass storage and this system will continue at DESY-IT for the LHC, photon-physics and others. Natural solution for DPHEP dataset
 - Changes “transparent” for user, relying on IT admin work



Key area 2: Software preservation & validation: sp-system

Write up: arXiv:1310.7814

- > Fairly early on, HERA experiments decided to try to migrate software for as long as possible rather than freezing the current environment
- > Pilot project of a system for software preservation and validation in 2010

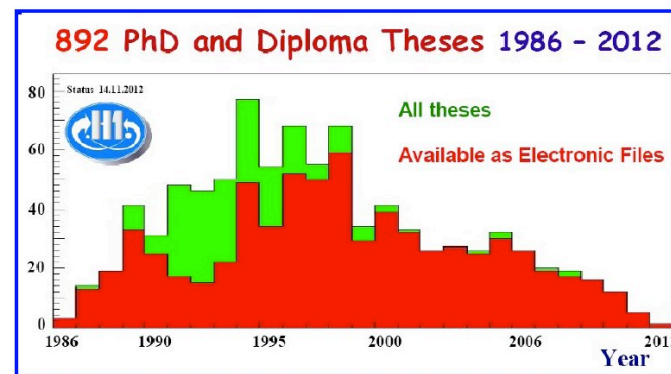
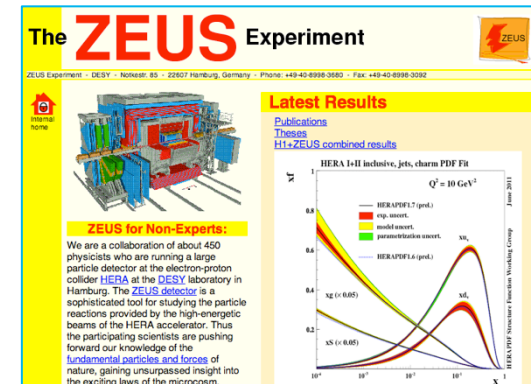
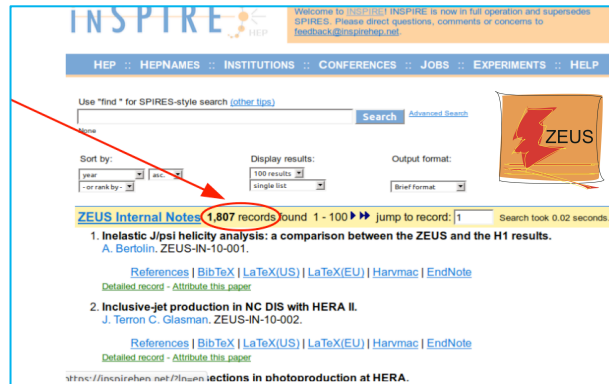


- > *Briefly:* The idea of the sp-system is to help perform migrations to newer software versions and environments, where transitions are performed often and validated by a comprehensive set of tests provided by the expts
 - The output of such a system is a **recipe** for deployment on (future) external resource(s)
 - Future analysis resources maybe local batch farm, grid, cloud, whatever
 - The idea is **not** to run analysis within the system itself!
- > Due to available resources and changes in personnel, implementation at DESY is still not in production mode
 - Project is rather ambitious and has taken longer than anticipated: definition of tests essentially done, but still requires much work to be done on the validation side



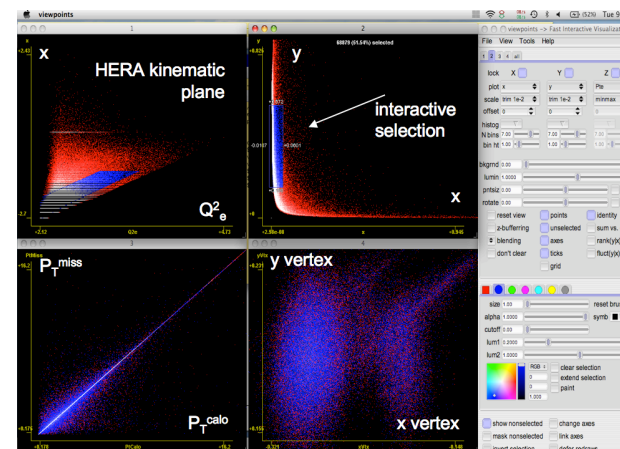
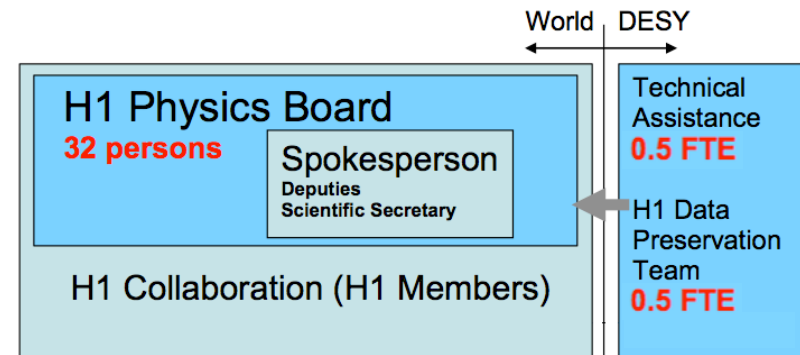
Key area 3: Documentation

- Successful collaboration between **INSPIRE**, the experiments and the DESY Library
- **Digital documentation** such as web-pages revised, reduced and streamlined for future use
- Lots of effort done sorting the vast amount of **non-digital documentation**
 - Many new (re-)discoveries along the way!
- Work done by key people with the *right expertise* and *experience* for the job



Key area 4: Governance, open access and outreach

- > H1 collaboration moved to a **new management model** in July 2012
 - Formation of *H1 Physics Board*, to replace Collaboration Board (institute based)
 - Future author list policies also set down in new constitution approved by collaboration
- > ZEUS and HERMES management teams retain same model as before, but similarly to H1 the collaborating institute layer is now removed
 - Remaining physics ZEUS working groups consolidated to a single physics group
- > **Open access** still to be considered and/or defined by the HERA experiments
- > **Outreach** is a great idea, but was not possible without dedicated resources
 - Already dropped in 2011 table shown earlier
 - Ideas existed, but nothing concrete came of it



One lesson already: HERMES



- > HERMES financial support officially ended December 31st, 2012
- > That year was busy time to try to finish off as much as possible, in terms of physics and data preservation
- > Still a wealth of interesting physics in their data!
- > Hardware turn-off and transfer to DESY-IT central services completed ✓
- > Validation project within `sp-system` not really implemented ✗
- > Current situation has no dedicated manpower at DESY for any HERMES activities, including data preservation
- > The same will apply to H1 and ZEUS, at least for DP, at the end of 2014



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)

- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)

- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!

- There is a **great reduction** in person power (and available expert knowledge) as well as funding as soon data taking stops. Budgets become much tighter towards the end, competing with other projects that are just beginning



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**
- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)
- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!
- There is a **great reduction** in person power (and available expert knowledge) as well as funding as soon data taking stops. Budgets become much tighter towards the end, competing with other projects that are just beginning
 - **Don't start too late**, projects should be well in place before data taking ends



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)

- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!

- There is a **great reduction** in person power (and available expert knowledge) as well as funding as soon data taking stops. Budgets become much tighter towards the end, competing with other projects that are just beginning
 - **Don't start too late**, projects should be well in place before data taking ends
 - **Don't underestimate the required person-power**: for funding or practical reasons



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)

- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!

- There is a **great reduction** in person power (and available expert knowledge) as well as funding as soon data taking stops. Budgets become much tighter towards the end, competing with other projects that are just beginning
 - **Don't start too late**, projects should be well in place before data taking ends
 - **Don't underestimate the required person-power**: for funding or practical reasons
 - **Dedicated manpower** is needed, people working on this part time or in spare time is not enough: such initiatives cannot “run for free”



Conclusions: Some lessons learned from DP @ DESY

- The **physics output** tail seen by LEP also rings true for the experiments at HERA, where there is much physics output in the years after data taking stopped
- In addition, the final data for preservation is not ready immediately after data taking
- Data volume, when the final data are available, **may not be such a decisive issue**

- Getting all of the data for preservation **to the same level** is quite some work but absolutely necessary: OS, software version (in house and external), calibrations, methodologies..
- This should be started **as soon as possible** (H1-HERA1 took 3 years to get 3 months of work)

- Best to **avoid dedicated material solutions**: use what's currently available. And who knows what this may look like in the future, so keep as generic and flexible as possible!

- There is a **great reduction** in person power (and available expert knowledge) as well as funding as soon data taking stops. Budgets become much tighter towards the end, competing with other projects that are just beginning
 - **Don't start too late**, projects should be well in place before data taking ends
 - **Don't underestimate the required person-power**: for funding or practical reasons
 - **Dedicated manpower** is needed, people working on this part time or in spare time is not enough: such initiatives cannot “run for free”
 - Losing the best people for the best roles is almost inevitable and finding support for unfinished things is extremely difficult. Difficult to capture the best candidates without providing a **long term perspective**

