

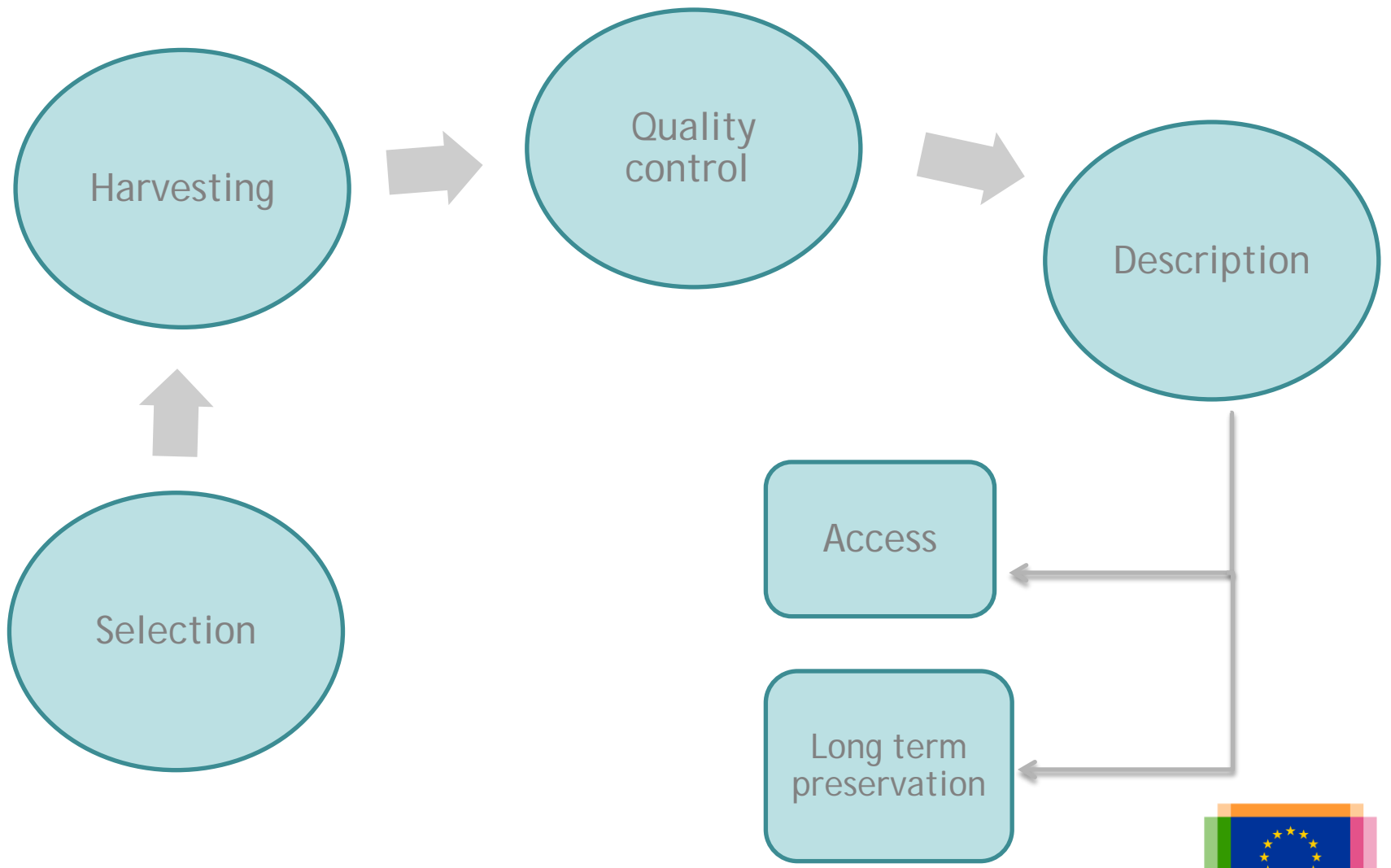


# Creating a web archive for the EU institutions' websites Achievements and challenges




Els Breedstraet  
Coordinator web preservation  
Publications Office of the EU

IIPC Web Archiving Conference, London, 16 June 2016

# Creating the archive






# Selection

	Internal: Publications Office / Historical Archives of the EU
	<ul style="list-style-type: none"><li>• Seed list of ca 70 europa.eu URLs</li><li>• Covers mainly sites of EU institutions, bodies and agencies</li></ul>
	<ul style="list-style-type: none"><li>• Selection policy<ul style="list-style-type: none"><li>• To be limited to europa.eu domain and subdomains?</li><li>• Which seeds to add? Which selection criteria?</li><li>• Who decides to add or not?</li></ul>=&gt; need for policies + curator and/or selection board</li><li>• Seed management:<ul style="list-style-type: none"><li>• How to propose new seeds =&gt; need for selection tool</li><li>• How to keep list up to date?</li></ul></li></ul>






# Harvesting

	Outsourced (Internet Memory Research)
	<ul style="list-style-type: none"><li>• Quarterly crawling of seed list</li><li>• Using Heritrix</li><li>• Runs successful since end of 2013</li><li>• Till mid 2016: ARC, afterwards WARC</li></ul>
	<ul style="list-style-type: none"><li>• Dynamic content</li><li>• Social media</li><li>• Multilingual harvesting</li><li>• Ad hoc demands/campaigns<ul style="list-style-type: none"><li>• How to decide to do or not?</li><li>• Planning and budgeting</li></ul></li></ul>






# Quality control

	<p>Automatic: outsourced Manual: website owners, coordinated by the Publications Office</p>
	<ul style="list-style-type: none"> <li>• Automatic checks + manual / visual control</li> <li>• Network of quality controllers</li> <li>• Structured follow up of feedback =&gt; long term improvement harvests</li> </ul>
	<ul style="list-style-type: none"> <li>• Automatic checks: more? Reporting?</li> <li>• Manual control             <ul style="list-style-type: none"> <li>• Need for quality control tool</li> <li>• Need for crawl reports</li> </ul> </li> <li>• Awareness raising -&gt; webmasters to be encouraged to make "well preservable" sites</li> </ul>






## Description / metadata

	Automatic (up to now)
	<ul style="list-style-type: none"><li>• Automatic capturing of descriptive and technical MD at the source</li><li>• Proposal for minimal set of MD (descriptive, technical and provenance) ready</li></ul>
	<ul style="list-style-type: none"><li>• Awareness raising -&gt; webmasters to be encouraged to add more and/or more structured MD at the source</li><li>• Provenance MD</li><li>• Post-crawl manual description?</li><li>• How to exploit MD to improve access?</li></ul>






# Access

	Outsourced (Internet Memory Research)
	<ul style="list-style-type: none"><li>• Fully open archive: <a href="http://collections.internetmemory.org/haeu">http://collections.internetmemory.org/haeu</a></li><li>• Basic search and browse options</li></ul>
	<ul style="list-style-type: none"><li>• Audience / use cases?</li><li>• Usage statistics</li><li>• Branding / customisation</li><li>• Data protection, copyright, etc.</li><li>• Multilingual access</li><li>• Promotion and communication</li><li>• Centralised access to EU related web archives? Cooperation?</li></ul>



# Long term preservation

	Internal: digital repository at Publications Office
	<ul style="list-style-type: none"><li>• First tests started with ingestion in internal trusted long term digital preservation repository</li></ul>
	<ul style="list-style-type: none"><li>• Adaptation information architecture</li><li>• Provenance metadata</li><li>• Trustworthy Digital Repository Certification (ISO16363)</li></ul>





# Conclusions

- Solid basis for an EU institutional web archive is in place, accessible to everyone
- Many challenges, a lot of work ahead
- Many dreams and ideas for the future
  - One step at the time
  - Cooperation?



# More information



Have a look at our web archive:  
<http://collections.internetmemory.org/haeu>

Questions or suggestions? Contact us:  
[op-web-preservation@publications.europa.eu](mailto:op-web-preservation@publications.europa.eu)



**Thanks for your attention!**

**Questions?**

