



SCAPE

Ambientes de preservação escaláveis

4º Seminário sobre informação na Internet // Preservação digital
21 Nov. 2012

Miguel Ferreira, PhD.
Director técnico // KEEP SOLUTIONS
mferreira@keep.pt

O que é o SCAPE?

Projecto de I&D em preservação digital

Visa o **planeamento** e **execução** de processos de preservação em **grandes colecções** de dados

Resultados esperados

Definição de **cenários** de preservação (para guiar os desenvolvimentos e validar os resultados do projecto)

Desenvolvimento de **software/ferramentas** open-source de auxílio à preservação digital

Criação de **workflows** de preservação

Edificação de um **infra-estrutura** de suporte à preservação (escalável)

Criação de **Guias de boas-práticas** em preservação (migração de repositórios e preservação de dados científicos)

O SCAPE é uma continuação de um anterior projecto chamado **PLANETS**

Financiamento

7º Programa Quadro (instrumento europeu de **financiamento de investigação**)

Projecto Integrado (IP) // Tecnologias da Informação e Comunicação (TIC)

Call sobre “**Bibliotecas digitais e preservação digital**”

Duração

42 meses // 3.5 anos

Fev 2011 a Jul 2014

Orçamento

11.3 Milhões de euros // 14.3 Milhões de dollars // ~30 milhões de reais

76% financiado pela Comissão Europeia

Institutos & fundações



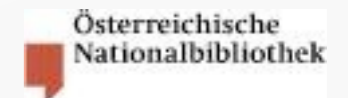
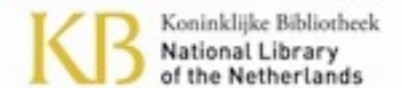
Universidades



Empresas



Bibliotecas



Detalhes

O SCAPE irá estender o estado da arte através da

Criação de uma **infra-estrutura** que permita a execução de **acções de preservação** em larga-escala (na ordem dos milhões de ficheiros e vários TB de informação)

Controlo de qualidade automático

Integração com **planeamento** de preservação suportado por **políticas** e **monitorização** automática do meio ambiente

Validação do projecto

O projecto será validado recorrendo a 3 bancos de ensaio

- ▶ 1. **Repositórios digitais** oriundos da comunidade das bibliotecas
- ▶ 2. **Conteúdos da Web** oriundos de 3 arquivos da Web
- ▶ 3. **Dados científicos** oriundos da comunidade científica

Open-source

Suportado quase exclusivamente por tecnologias open-source

Resultados serão publicados com licenças APACHE 2.0

Dados

Repositórios digitais

Conteúdo da Web

Dados científicos

Dados

Repositórios digitais

Conteúdo da Web

Dados científicos



Caracterização

(Identificação de formatos, extracção de propriedades)

Dados

Repositórios digitais

Conteúdo da Web

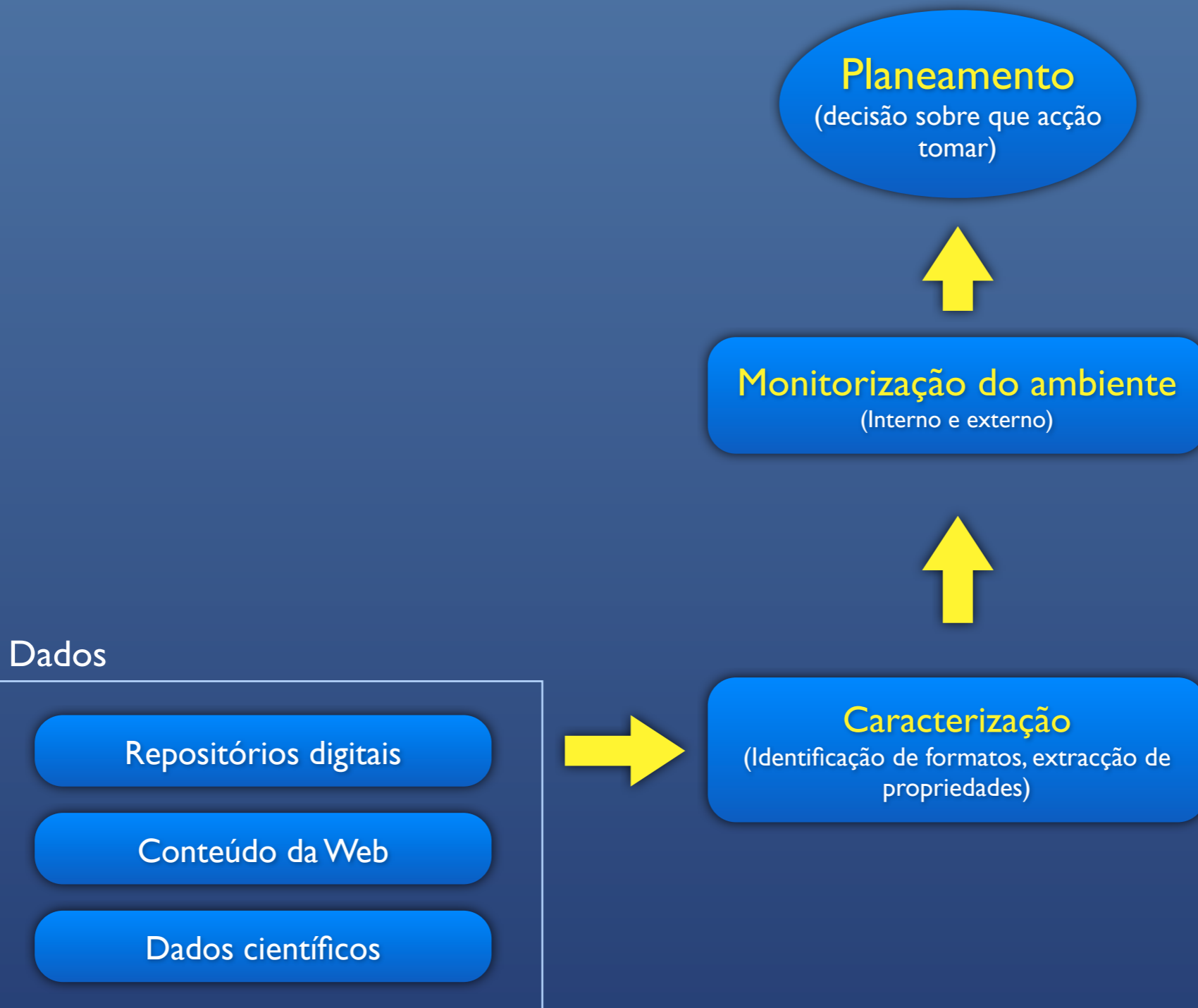
Dados científicos

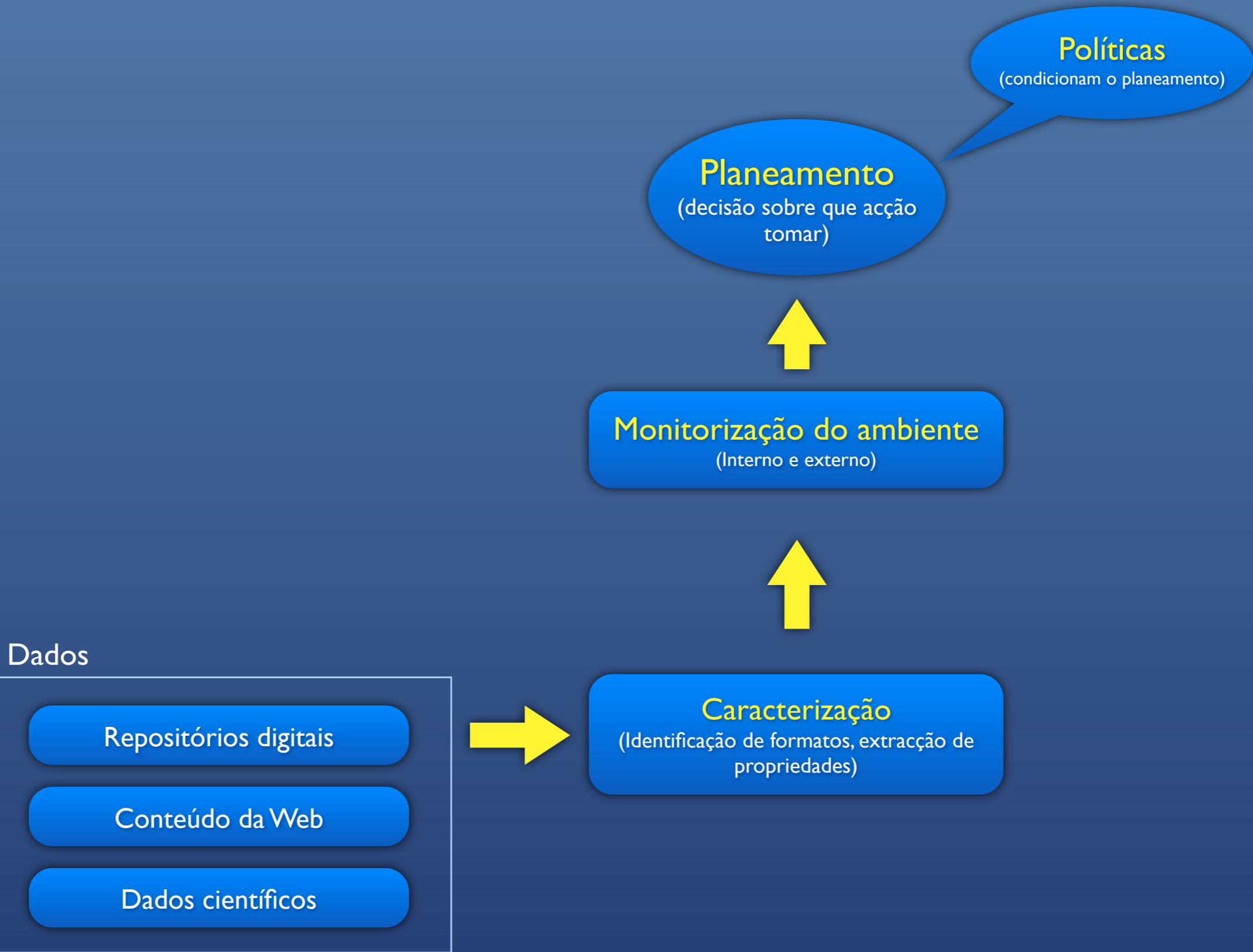


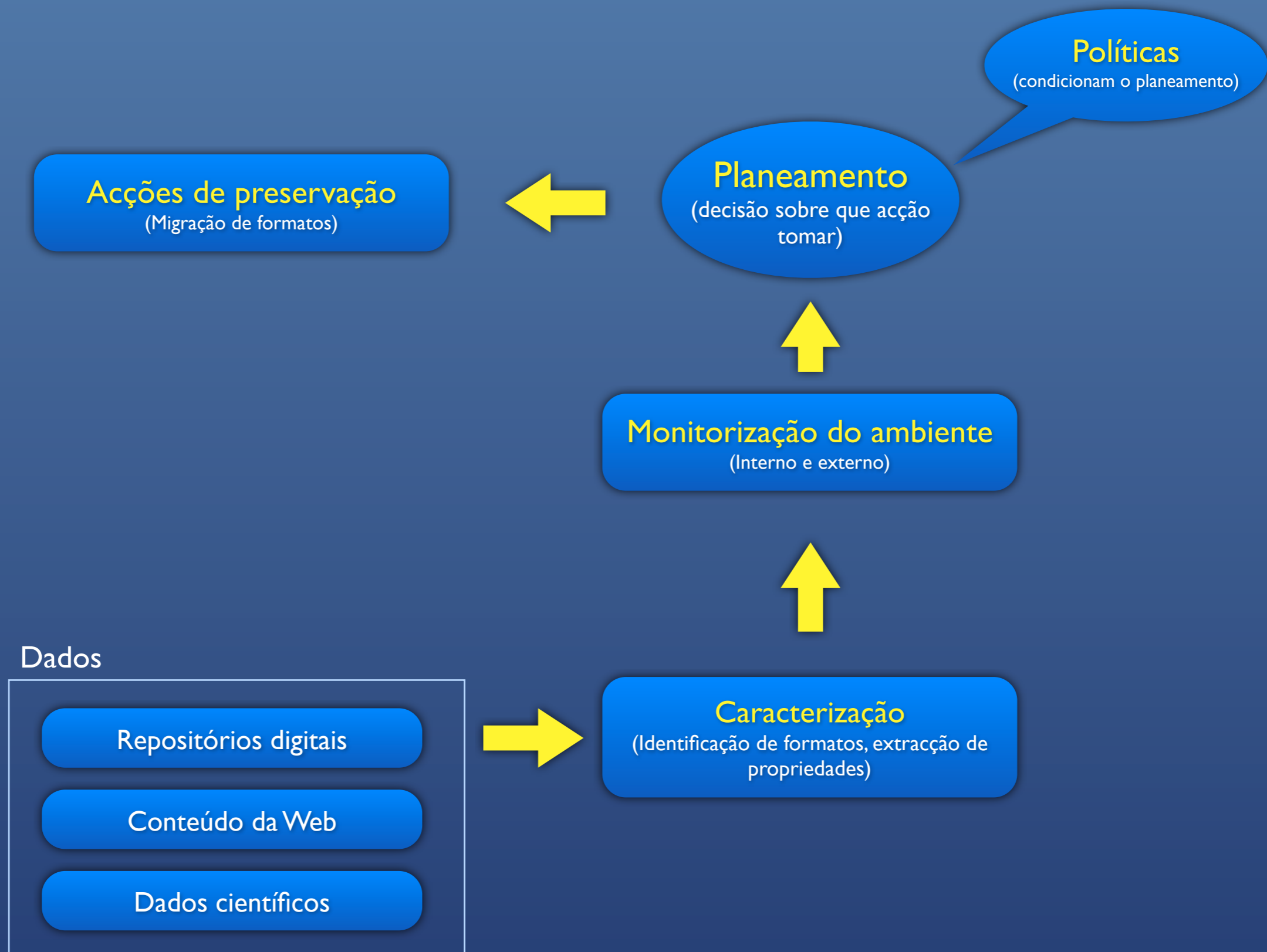
Monitorização do ambiente
(Interno e externo)

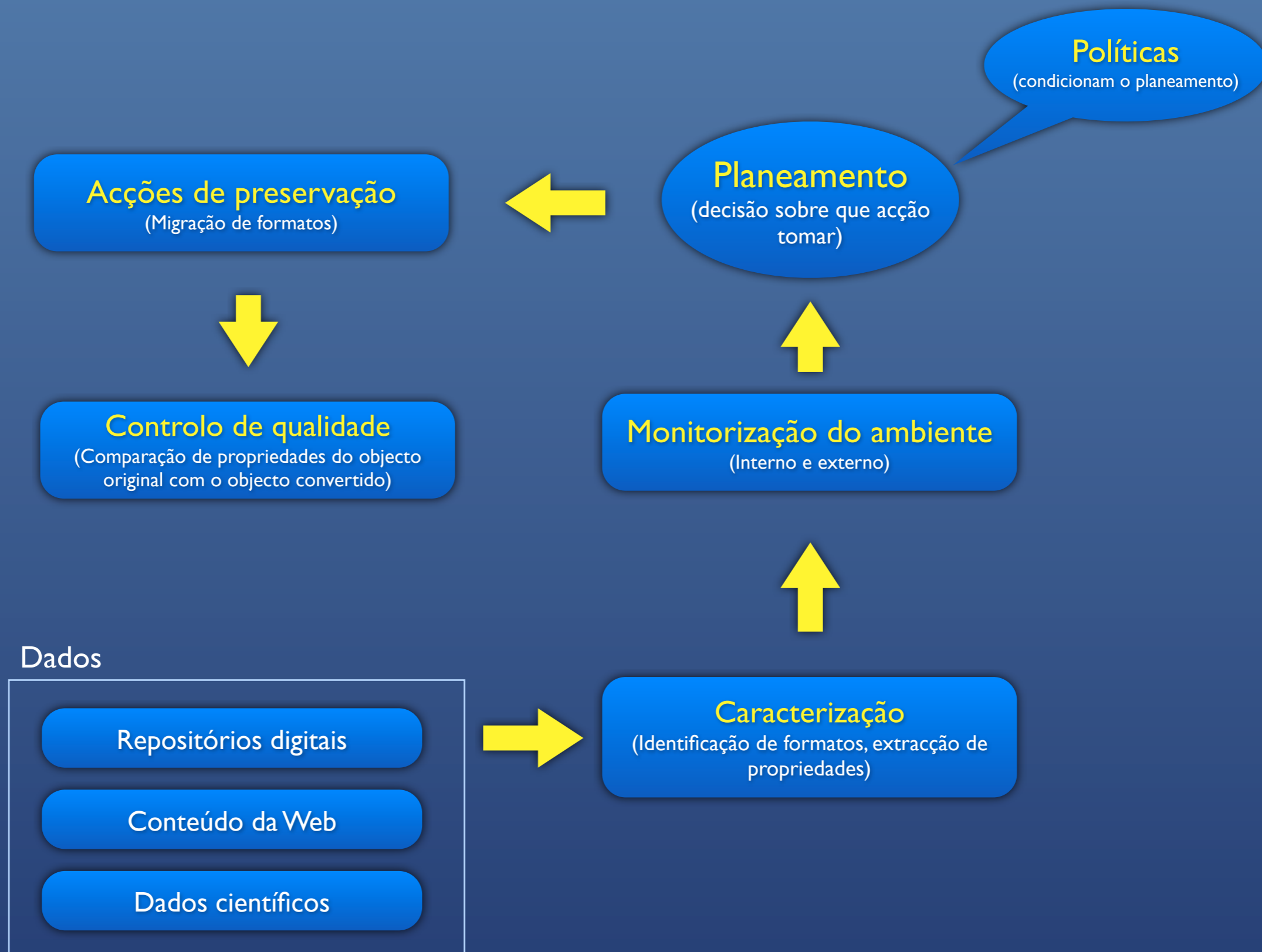


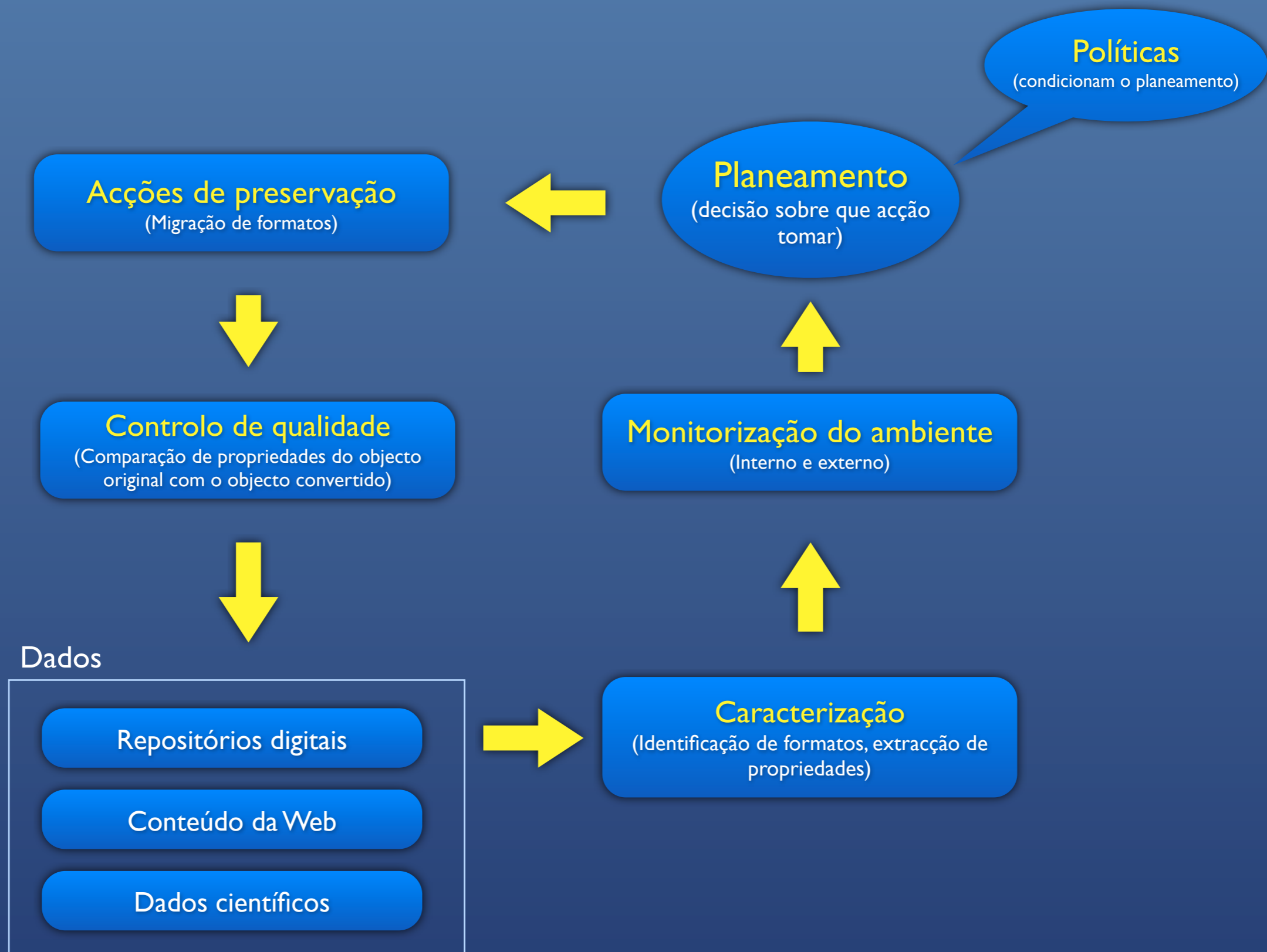
Caracterização
(Identificação de formatos, extracção de propriedades)

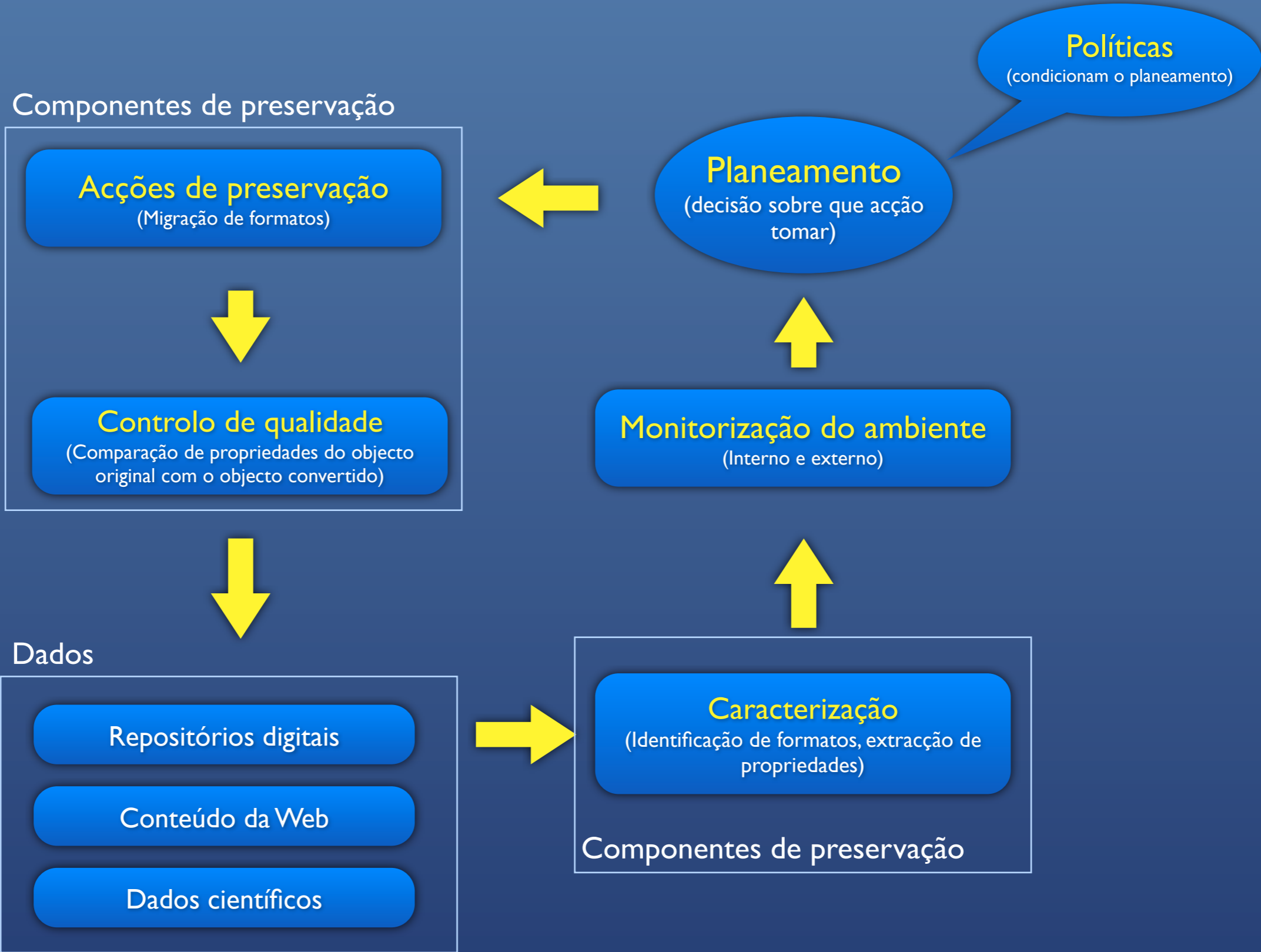


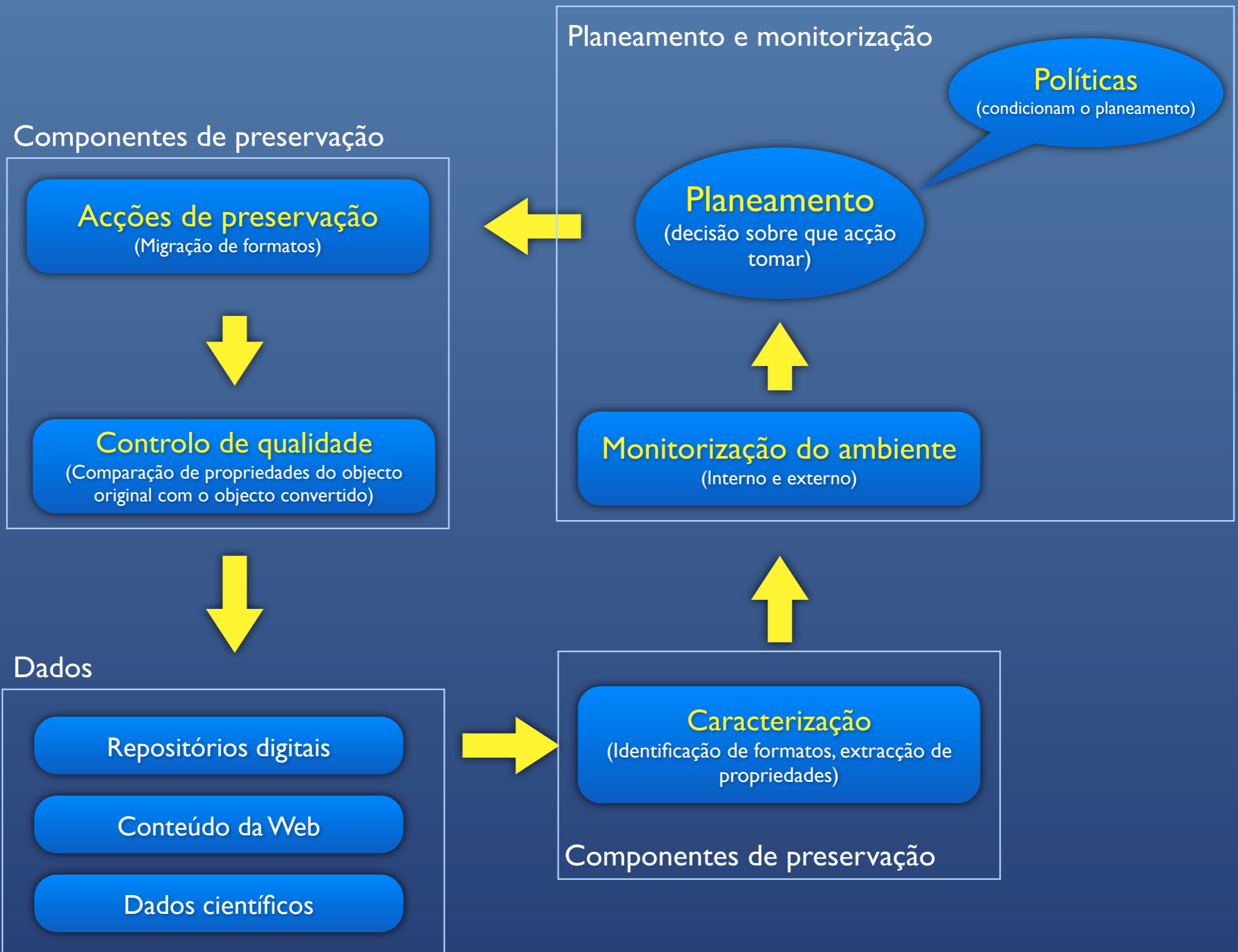




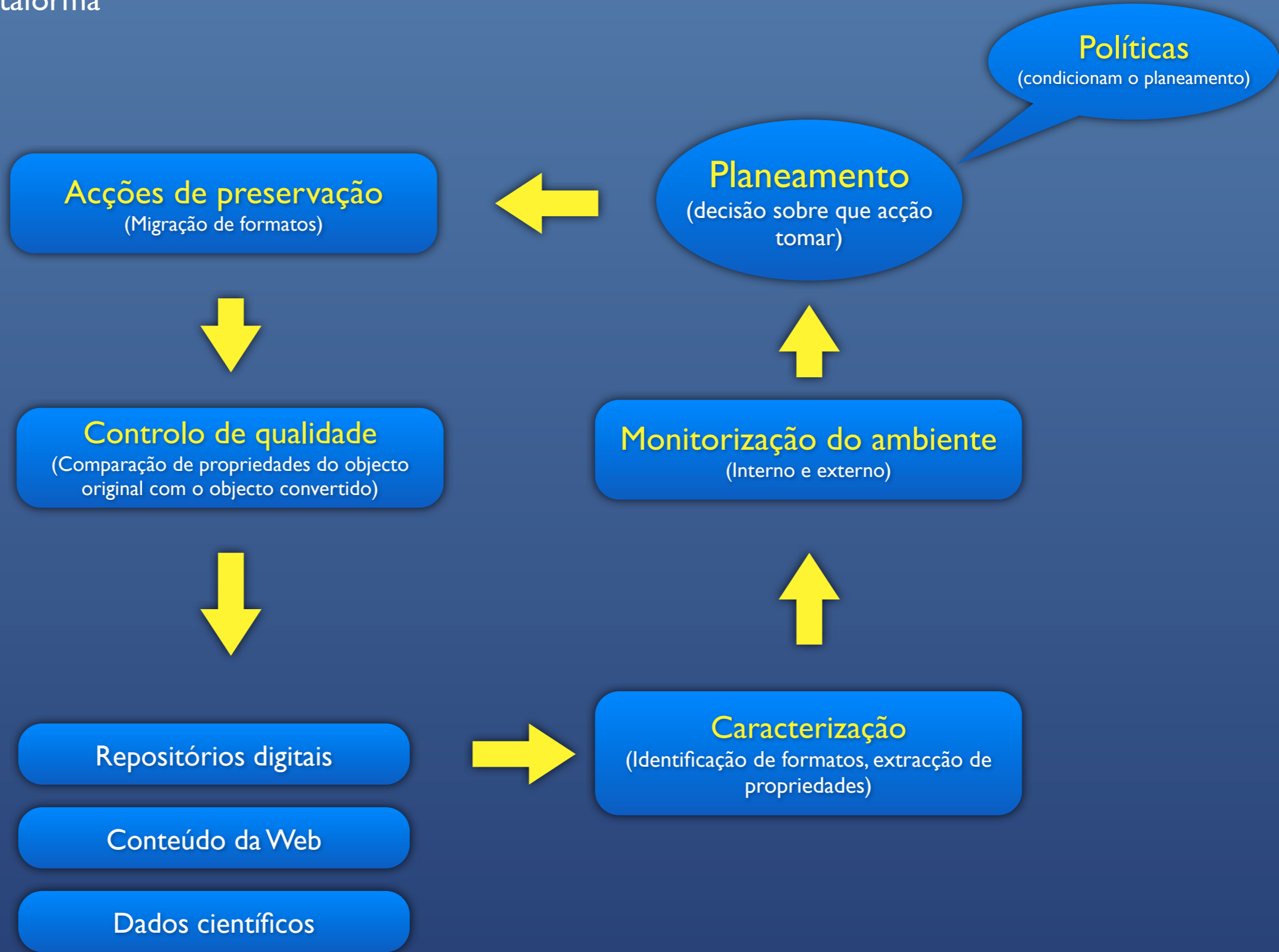


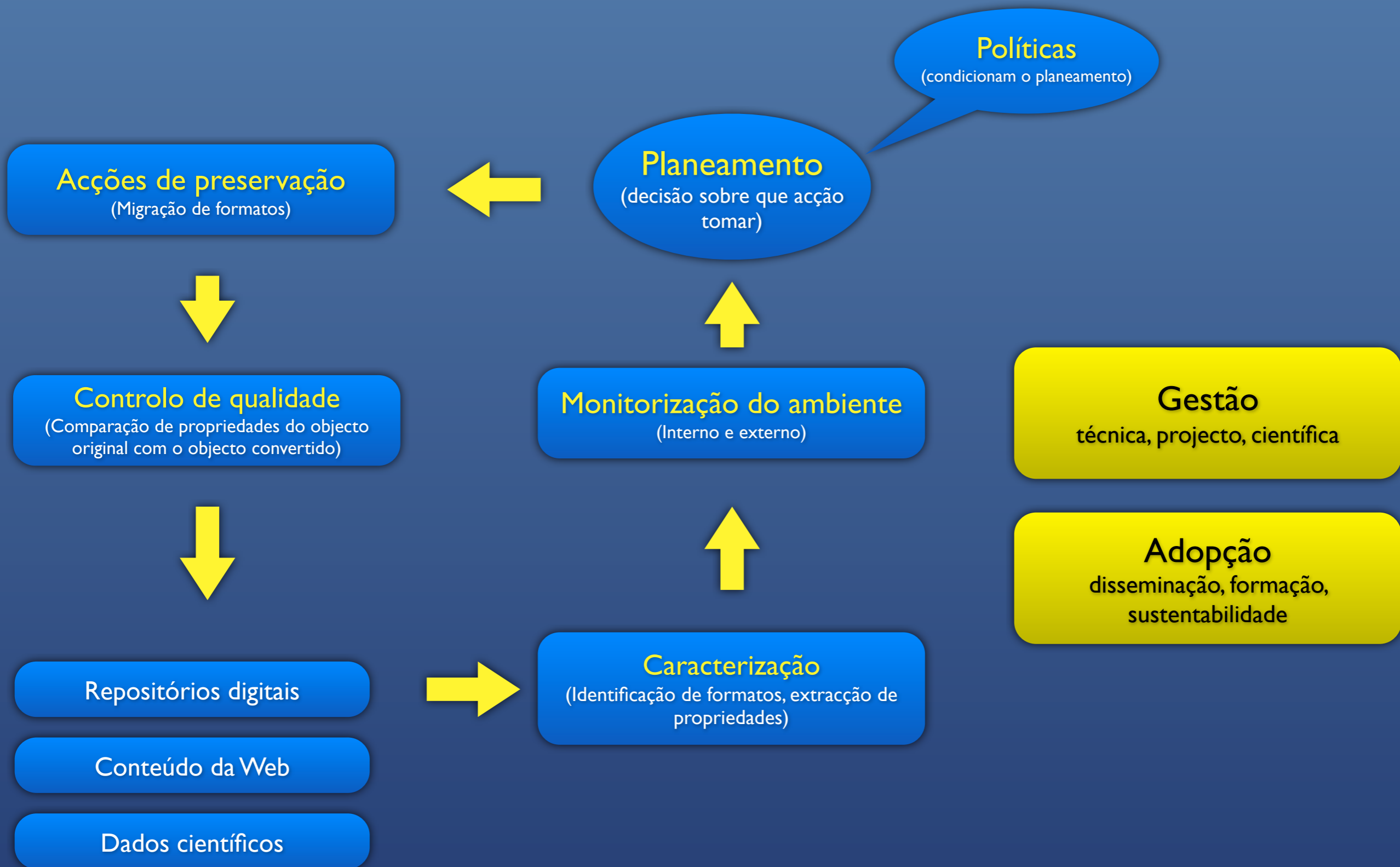






Plataforma





Dados e cenários de aplicação

3 bancos de ensaio

Repositórios digitais

Dados científicos

Arquivos da Web



Migração de TIFF para JPEG 2000

Redução de custos de armazenamento

80 TB // 2.2 M de páginas de jornais digitalizados do século XIX

Conversão de RAW para NeXus

Instrumentos de medição produzem dados em formato RAW. Pretende-se migrar todos esses dados para o formato NeXus baseado em XML por ser um standard internacional

Capacidade de processar ficheiros na ordem dos 100 GB



Controlo de qualidade em arquivos da Web

Conhecer o conteúdo do arquivo da Web para saber se é necessário agir para preservar

Melhorar os processos de recolha da Web de modo a garantir que **não ficaram partes importantes por recolher**

Componentes de preservação

Ferramentas de **caracterização**

Visam dar a conhecer os nossos dados

Identificação de formatos

Extracção de propriedades internas dos objectos

Nº de páginas, tamanho, largura, altura, esquema de cores, metadados, etc.

Exemplos: FITS, Tika, ffproble, Unix File, Droid, ...

Acções de preservação

Visam tornar os dados mais acessíveis

Migração de formatos

Emulação

Exemplos: ImageMagick, MS Sharepoint services, ffmpeg, gstreamer, ARC2WARC

Controlo de qualidade

Visam garantir que as acções de preservação não comprometeram a qualidade dos objectos

Comparam propriedades dos objectos antes e depois das acções de preservação terem sido aplicadas

Exemplos: Image Comparison, Matchbox, MarcAlizer VIPS, Jpylyzer

Digital Preservation Toolkit

<http://tinyurl.com/dptoolkit>

Planeamento e monitorização

Monitorização do meio ambiente

Visa detectar eventos internos e externos ao nosso repositório de informação que possam constituir riscos ou oportunidades relevantes para a preservação dessa informação

Suporte para diferentes fontes de informação

- ▶ directórios de formatos, directórios de software, repositórios, utilizadores, tendências, etc.

Permite exprimir riscos e oportunidades sob a forma de “questões”

Notifica os interessados sempre que o estado do “mundo” viola uma determinada condição

Planeamento de preservação

Avaliação das diferentes alternativas de preservação com vista a maximizar os objectivos da instituição

Plato - <http://www.ifs.tuwien.ac.at/dp/plato>



Políticas institucionais

Definição formal de políticas institucionais que poderão influenciar decisões ao nível do planeamento de preservação

- ▶ Todos os dados devem ser preservados. Nada é eliminado
- ▶ O custo de uma ferramenta de preservação não deve exceder X
- ▶ O custo do salário de um especialista em preservação não deve exceder Y
- ▶ Todos os conteúdos devem ter um checksum

Plataforma

Infra-estrutura onde todos os componentes são executados

Cluster **Hadoop**



Replicação de dados

Integração com repositórios

HDFS, HBASE, HDFS-Fedora



Três níveis de paralelização

Distribuição dos dados pelos vários nós do cluster (acesso mais rápido)

Divisão de **grandes ficheiros** para processamento paralelo

Paralelização de **algoritmos** e/ou do trabalho



Adopção

Disseminação

Sítio **Web**, Newsletter

Twitter

OPF **Blog**, Slideshare, Vimeo

Github + OPF wiki

Publicações científicas (22 artigos já publicados)

Materiais promocionais (autocolantes, **brochuras**, pins)

Formação

Várias **formações** previstas ao longo do projecto

A primeira será nos dias 6-7 em Guimarães, Portugal (INSCRIÇÕES ABERTAS)

Sustentabilidade

Evangelização através da **participação em eventos**

Integração de resultados em **aplicações comerciais**

Promoção do **uso** das aplicações

Mais info

Web Site

<http://www.scape-project.eu>

Código fonte

<https://github.com/openplanets/scape/>

Wiki de desenvolvimento

<http://wiki.opf-labs.org/display/SP/Home>

Workflows experimentais

<http://www.myexperiment.org/search?query=SCAPE&type=all&commit=Search>

Publicações

<http://www.scape-project.eu/category/publication>

Relatórios de projecto

<http://www.scape-project.eu/category/deliverable>

KEEP SOLUTIONS

University of Minho SPIN-OFF

www.keep.pt

Miguel Ferreira
mferreira@keep.pt

Acknowledgements

Miguel Arellano (IBICT)
Dr. Ross King (AIT)