

学习科研经历：

-2004 系统工程，大连理工大学
2004-2007 知识科学，北陆先端科学技术大学院大学
2007- IME，ASA 国际应用系统分析研究所

项目经验：

- ✓ Energy efficiency and Risk Management in Public Buildings (EnRiMa)
- ✓ New Energy Externalities Developments for Sustainability (NEEDS)
- ✓ Web-based Emission Trading System
- ✓ MCA-based Global Energy Assessment
- ✓ Intuitive Decision Analysis support System (Prototype)
- ✓ COE: Knowledge Management system, E-science Environments, Domain Knowledge Ontology Construction

大数据分析及其在 Web3.0 中的应用

Hongtao Ren

renh@iiasa.ac.at

International Institute for Applied Systems Analysis
A-2361 Laxenburg, Austria

ECUST, 8-12 Oct, 2012

大数据定义：

大数据是指需要即时处理的、数据集容量非常庞大的、结构非常复杂的数据。

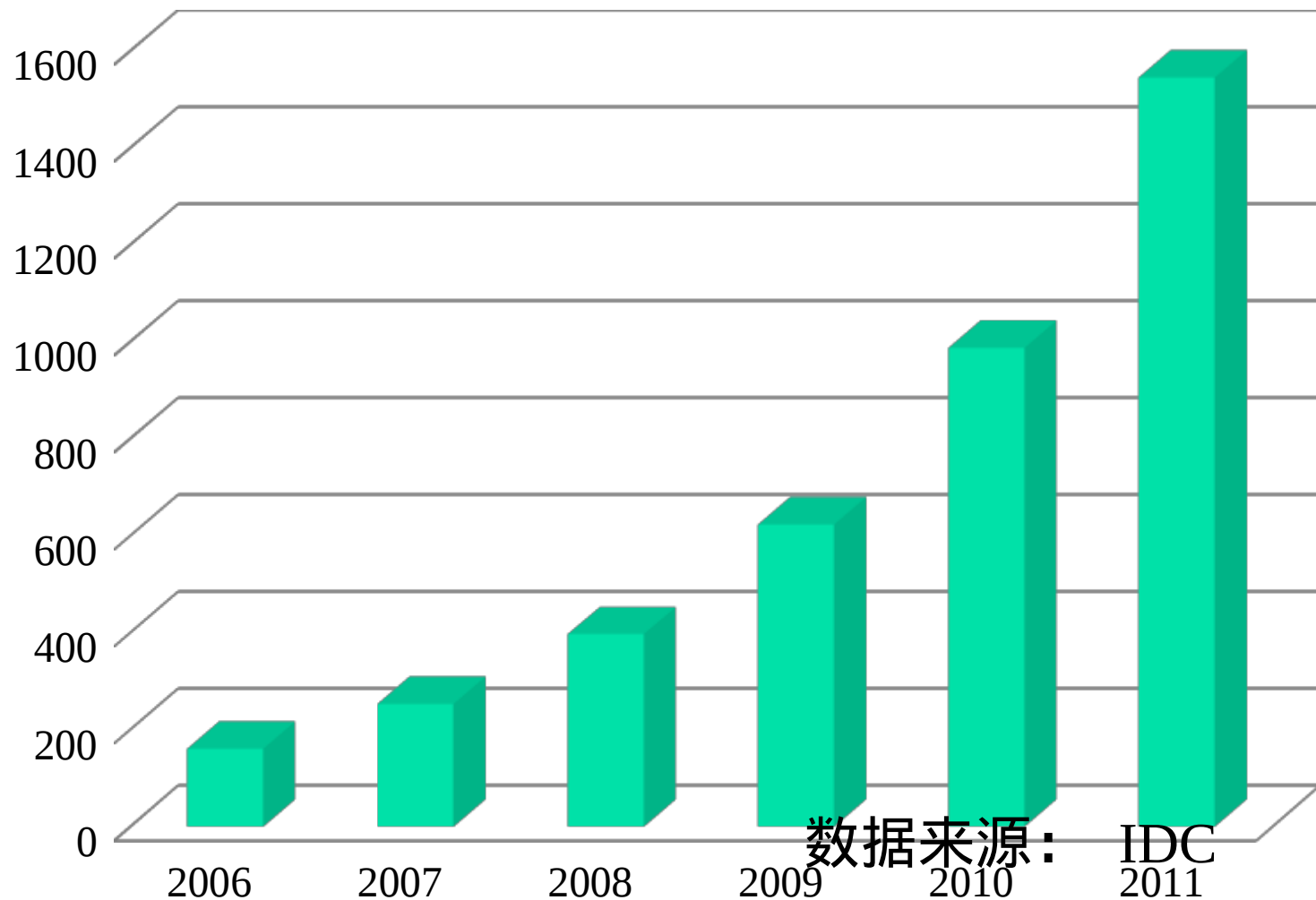
(volume, velocity, and variability)

数据处理包括：

- 数据捕捉
- 数据存储
- 数据搜索
- 数据分享
- 数据分析
- 数据可视化



每年存储的数据量 (10¹⁸ Bytes)



典型的大数据应用领域：

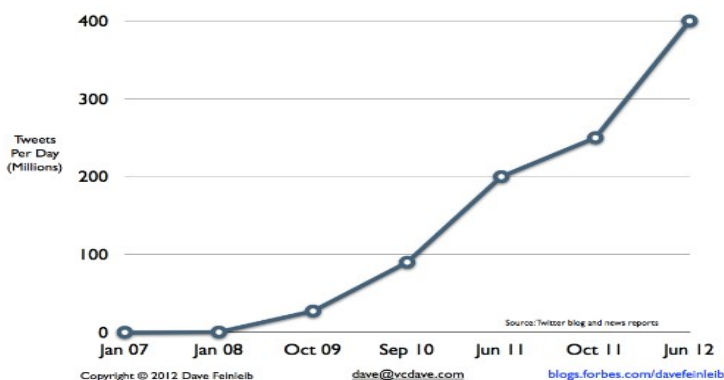
✓ 网络日志：

沃尔玛 2012 年平均每小时要处理大于 100 万客户的交易数据，数据容量约 2.5×10^6 G 字节的数据；

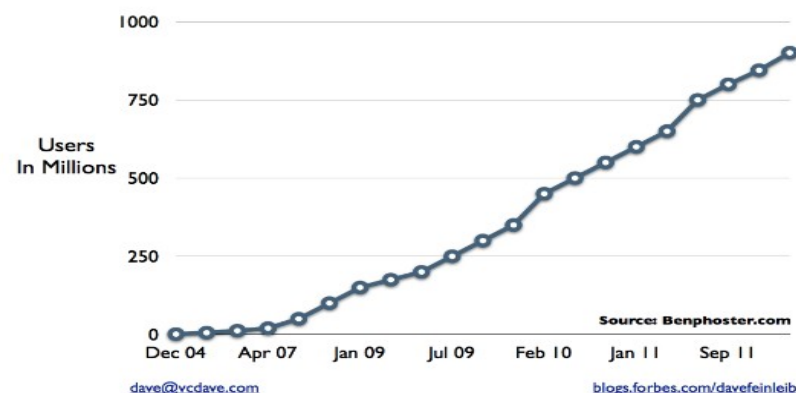
亿赞普每天要对 40×10^3 G 字节的日志进行数据挖掘、精确广告投送

✓ 社会网络：

Twitter: Tweets Per Day



Facebook User Growth



✓ 基因学，天文，地理信息

✓ 结构化数据：

- 1) 传统关系型数据库 Oracle , MySQL, postgresSQL, MSSQL
- 2) 二维表存储实体、关系模型
- 3) 1NF , 2NF , 3NF 来确保数据无冗余以及数据一致性

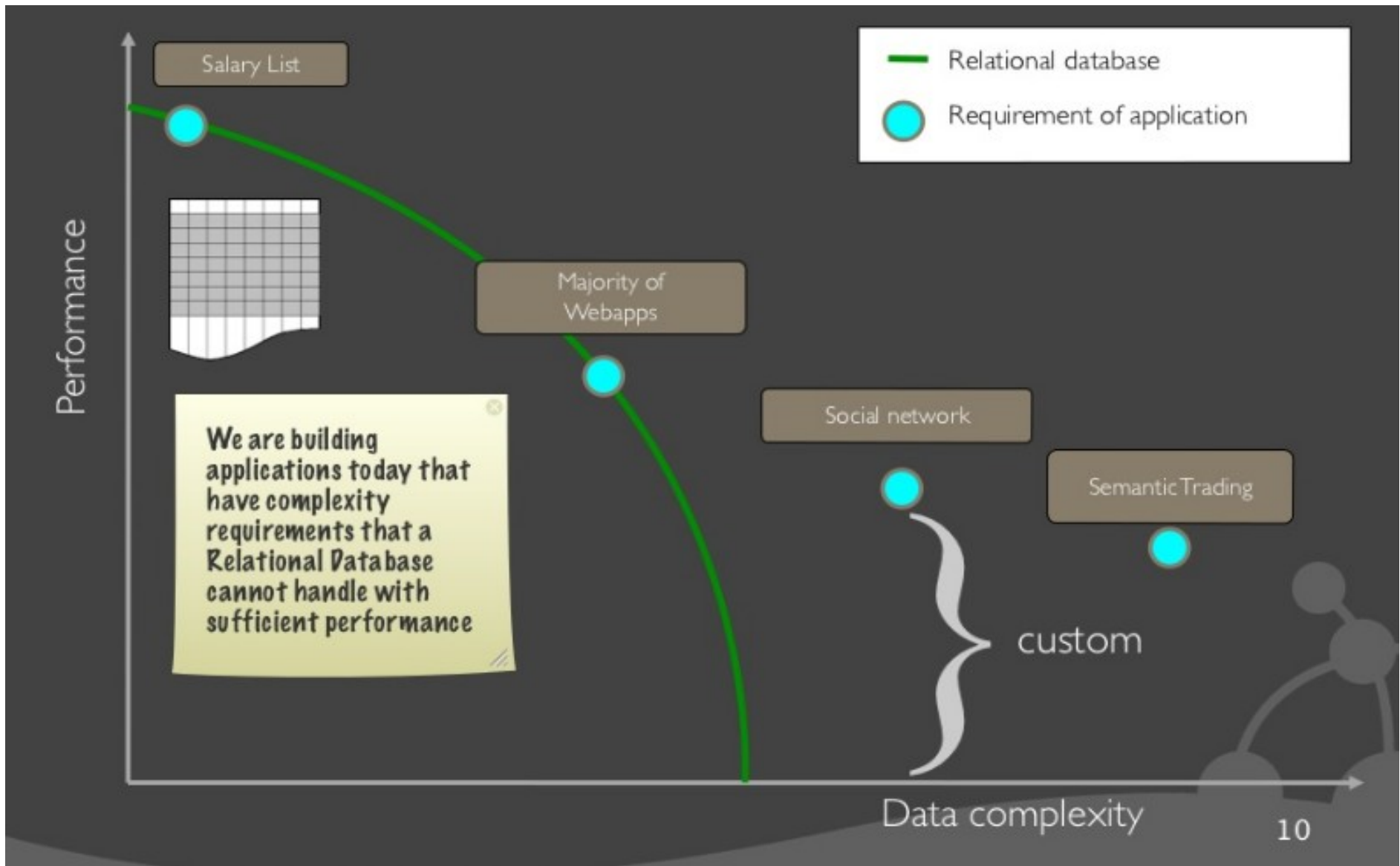
✓ 半结构化数据：

- 1) 面向对象数据库 (Versant 、 UNISQL)
- 2) XML , RSS , JSON

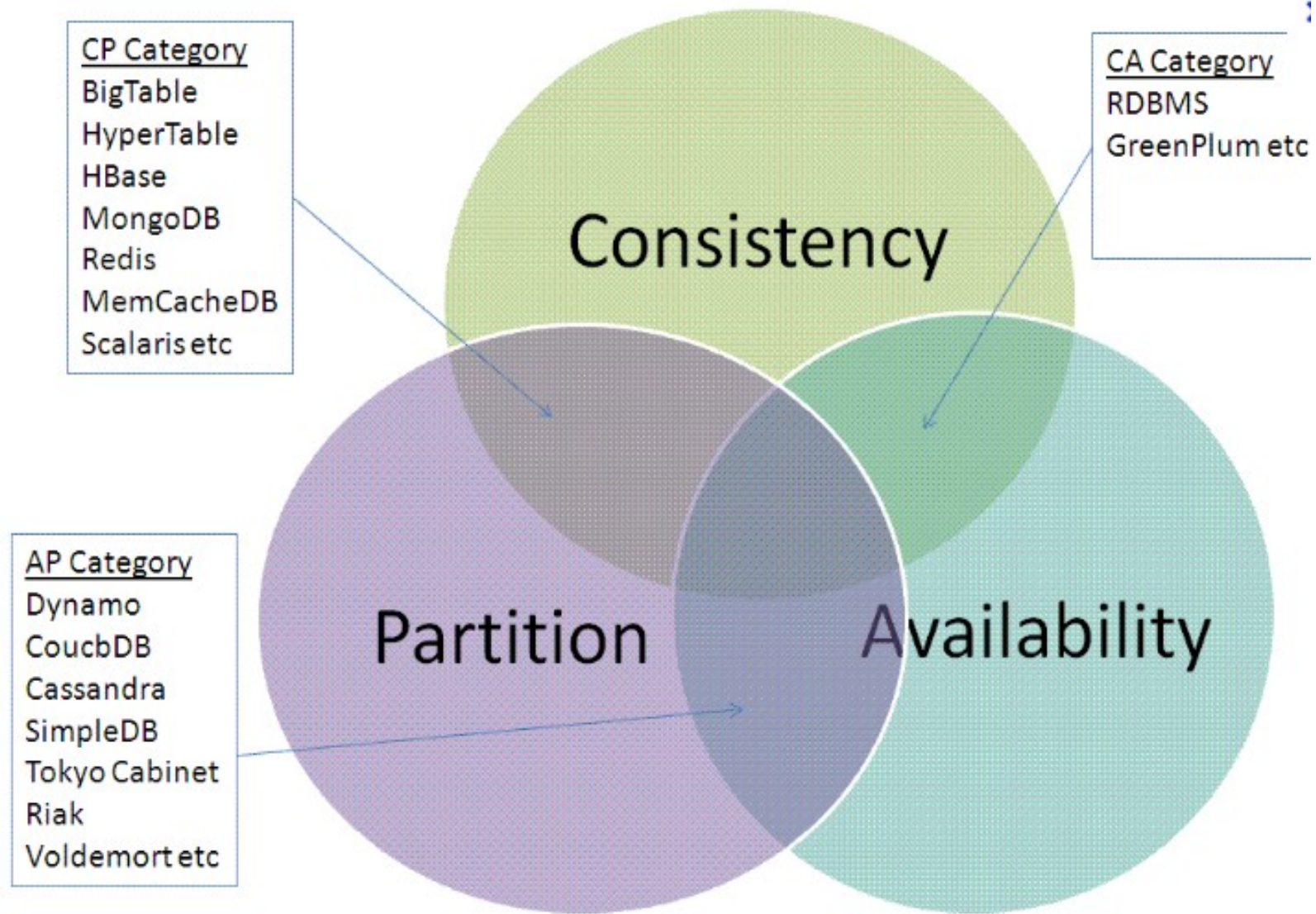
✓ 非结构化数据：

- 1) 网页
- 2) 语义网, RDF
- 3) Ontology
- 4) 标签、注释

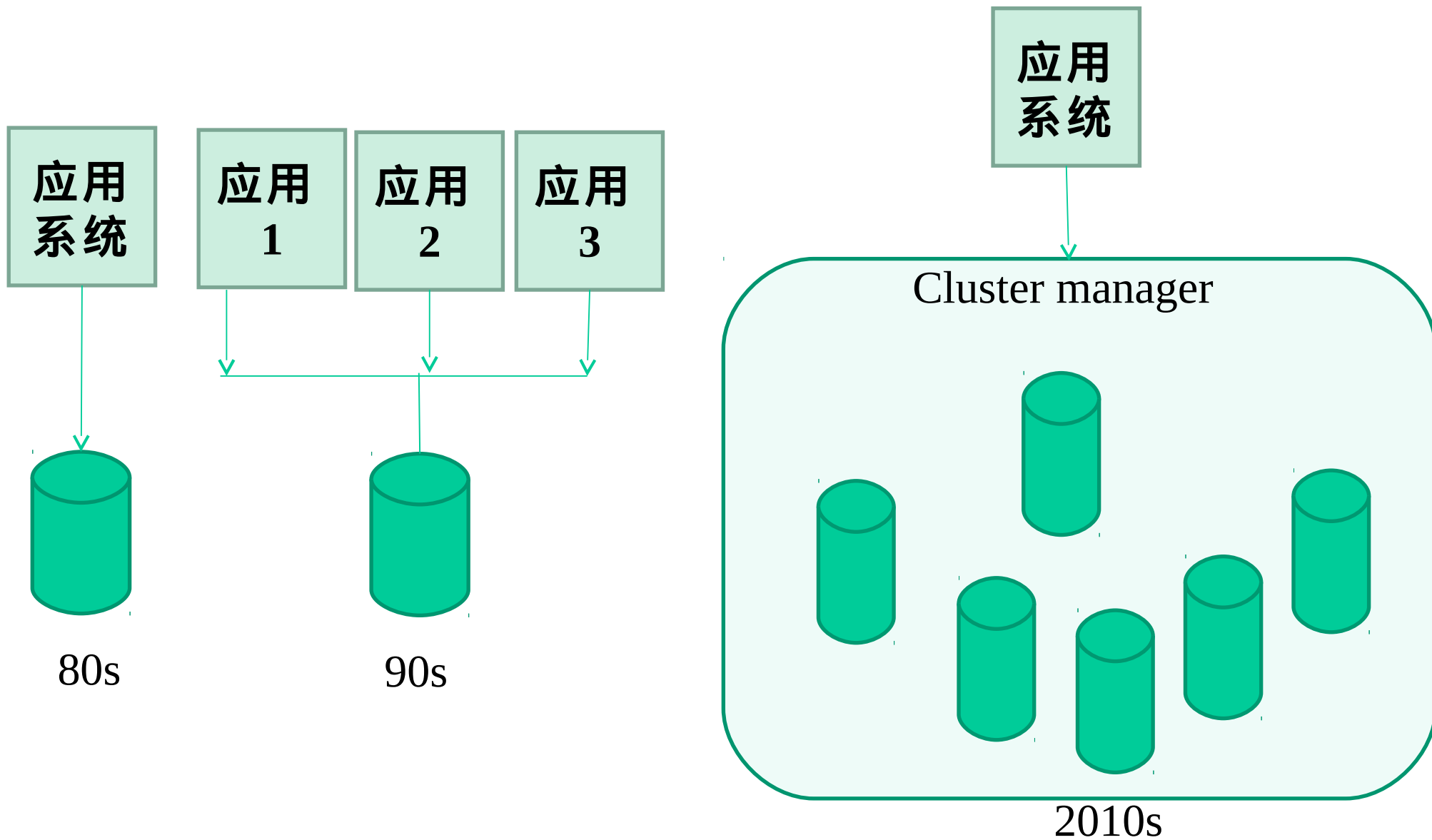
结构复杂性与查询效率

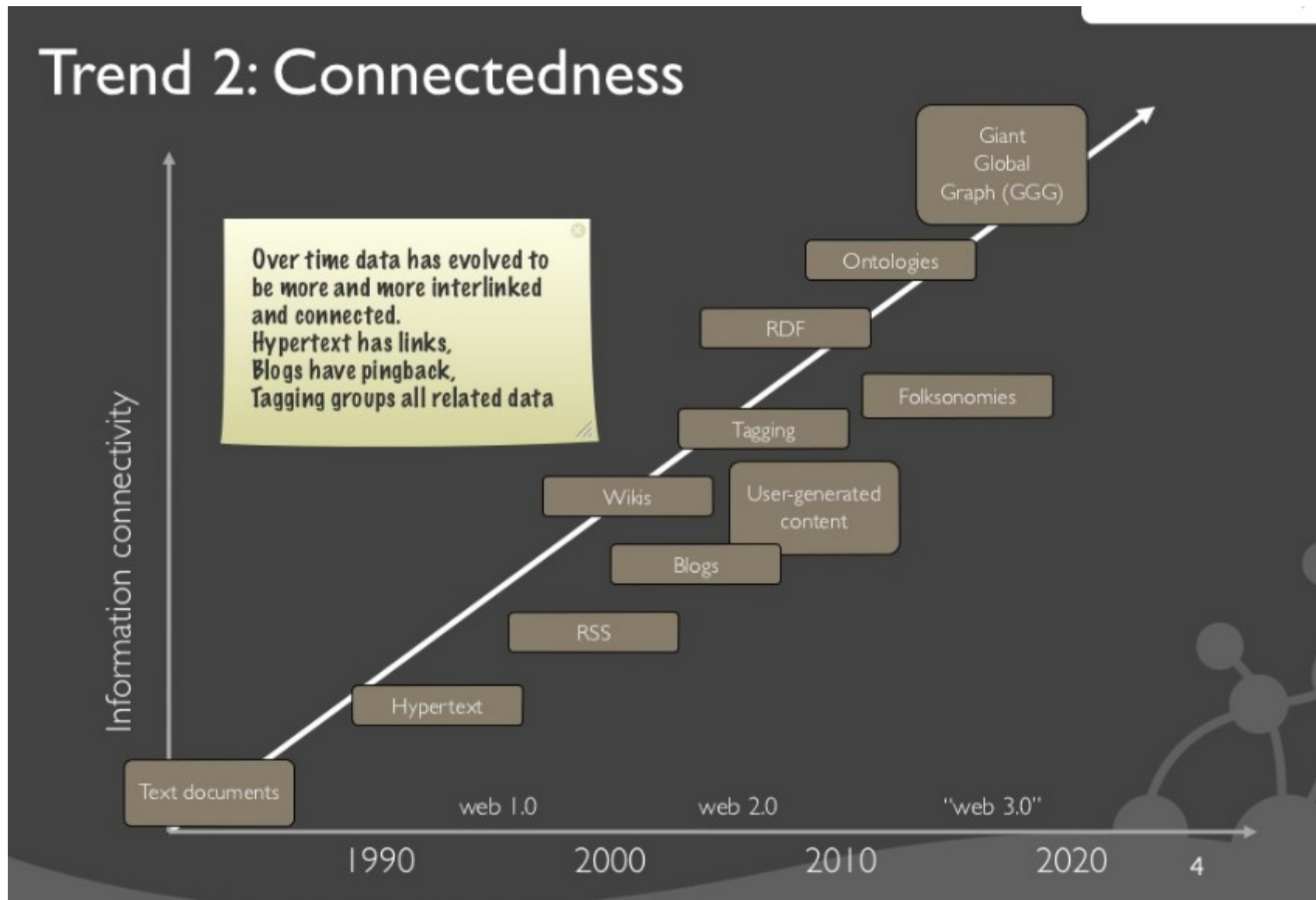


来源： Noe4j



CAP theorem





多准则问题：

Web 1.0

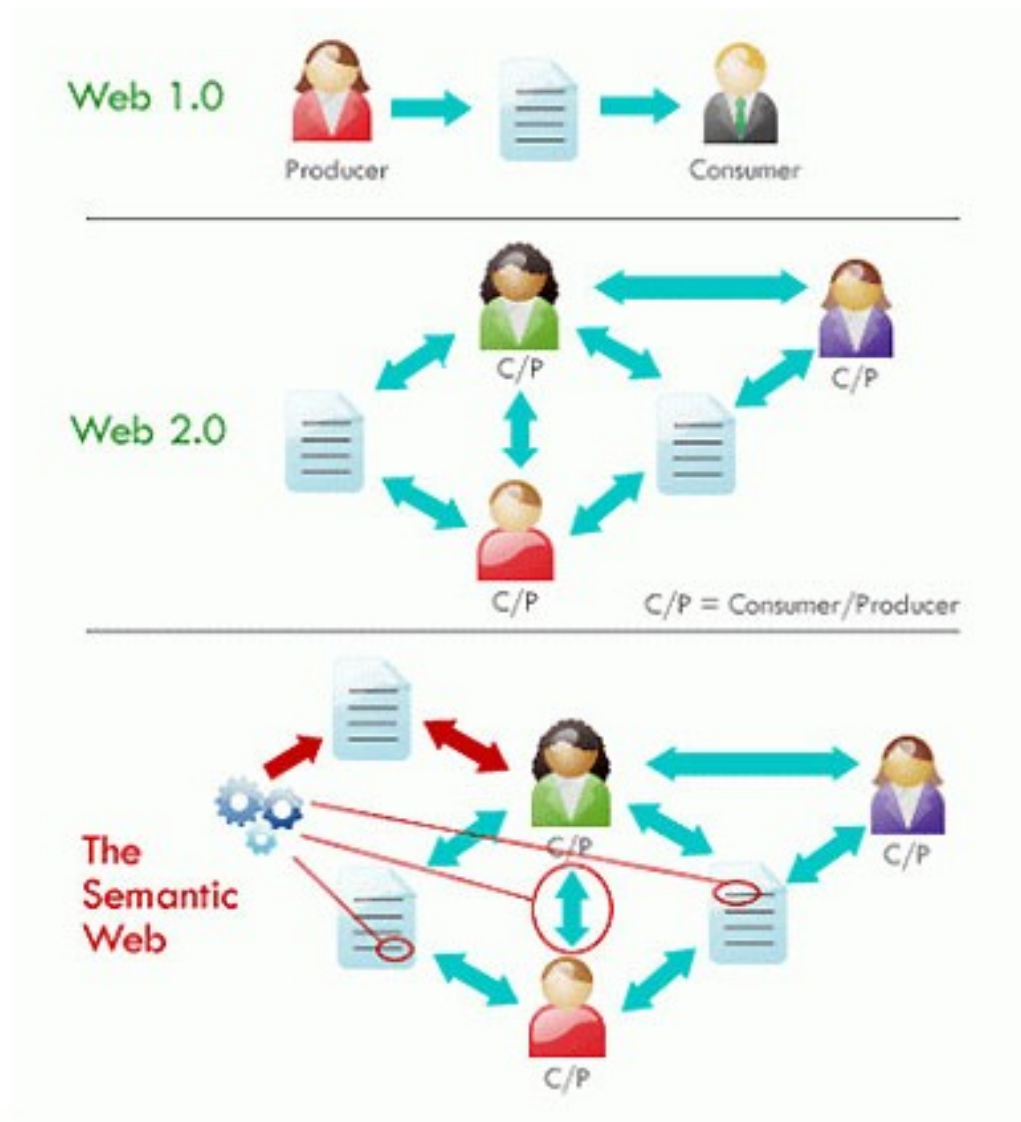
“the mostly read only web”
45 million global users (1996)
focused on companies
home pages
owning content
Britannica Online
HTML, portals
web forms
directories (taxonomy)
Netscape
pages views
advertising

Web 2.0

“the wildly read-write web”
1 billion+ global users (2006)
focused on communities
blogs
sharing content
Wikipedia
XML, RSS
web applications
tagging (“folksonomy”)
Google
cost per click
word of mouth

Web 3.0

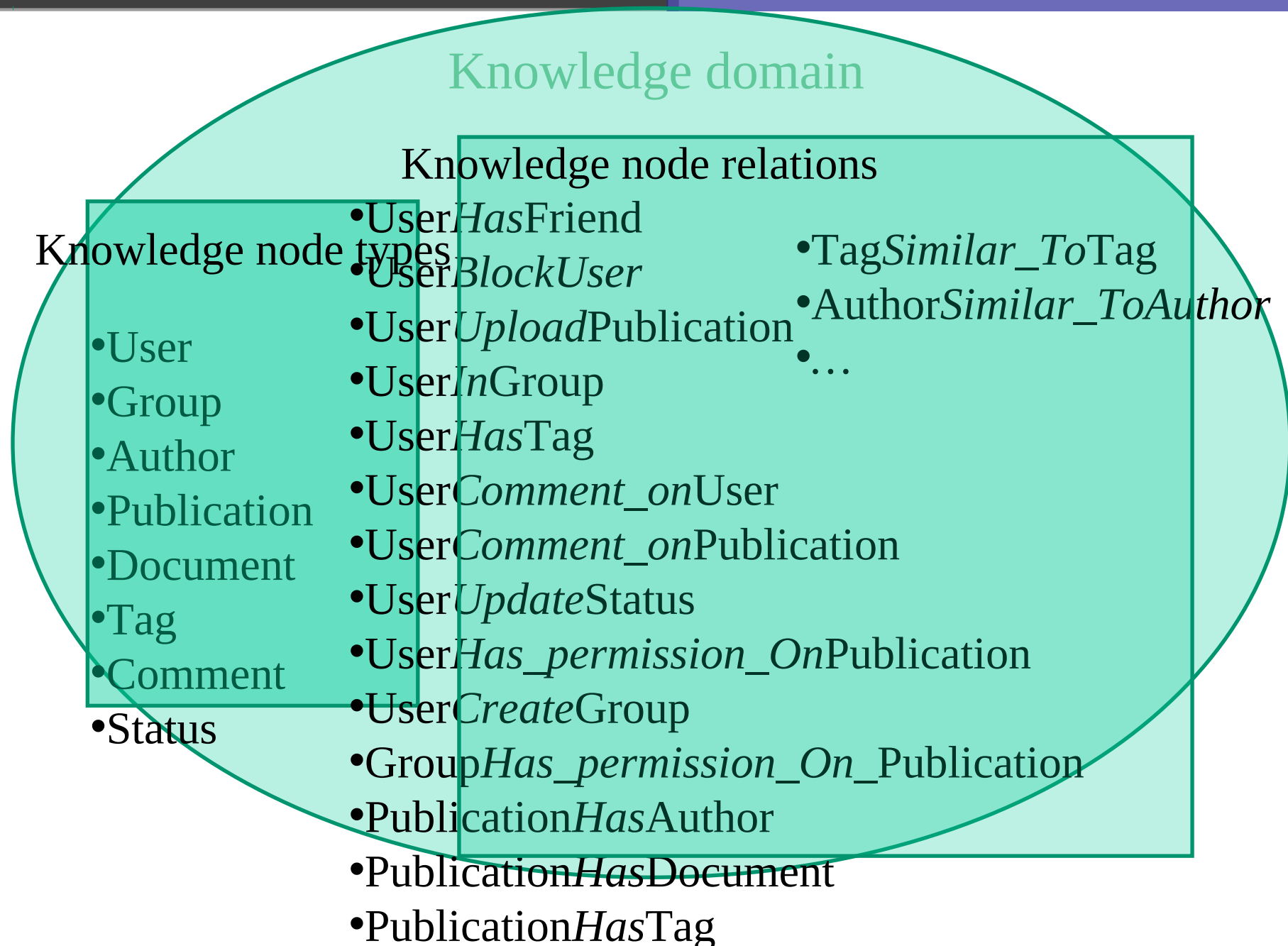
“the portable personal web”
focused on the individual
lifestream
consolidating dynamic content
the semantic web
widgets, drag & drop mashups
user behavior (“me-onomy”)
iGoogle, NetVibes
user engagement
advertainment



- A web where the context of the content is defined as data
- A web capable of reading and understanding content and context
- A web can better satisfy the requests of the people and machine
- A web capable of filtering the content that is of interest to the user

Challenges

- Connectedness: All different type of knowledge node are interlinked and connected
- Data size: For example, if we want to store 10 thousand publications, we may need store over 1 million related entities
- Semi-structure: Individualization of the knowledge node, the property of the knowledge node should be user definable
- Extracts implicit, potentially useful information from the data
- Knowledge visualization





- ✓ JDK7
- ✓ C++, PL solver
- ✓ Tomcat7
- ✓ Sun Grid Engine
- ✓ Springframework
- ✓ Axis2 Web services
- ✓ Jfreechart
- ✓ Hibernate
- ✓ Oracle10g
- ✓ Ajax

Layers:

- Java persistence layer
- Object-relational mapping abstraction layer
- Data access layer
- Services layer
- Process and integration layer
- Interfaces layer

Modules

- Problem
- Instance
- Analysis

Lessons from science-policy interactions:

- Much more demanding than pure science
- There is no golden key: modeling requires a combination of science, craft, art, experience
- Modeling for (interdisciplinary) knowledge integration and creation

A. Einstein:

Everything should be made as simple as possible, but not one bit simpler