

# Go with the dataflow

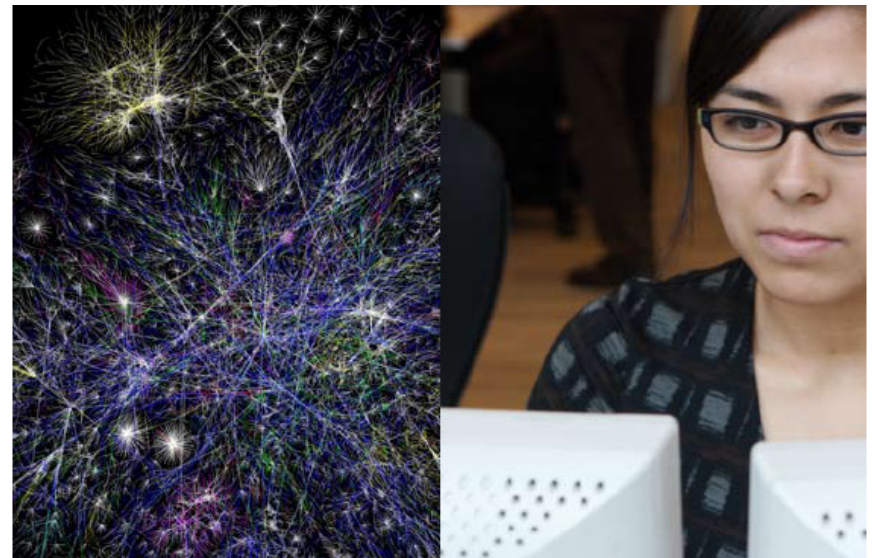
Presentation at the meeting of EUROSTAT  
Information Society Statistics Working Group



**Christiaan Holland**  
**Robbin te Velde**

Luxembourg, October 8<sup>th</sup> 2008

**dialogic**  
innovatie • interactie



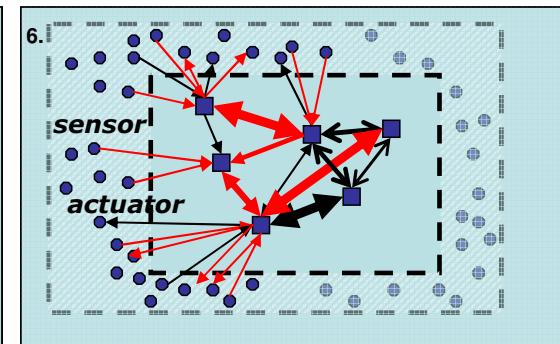
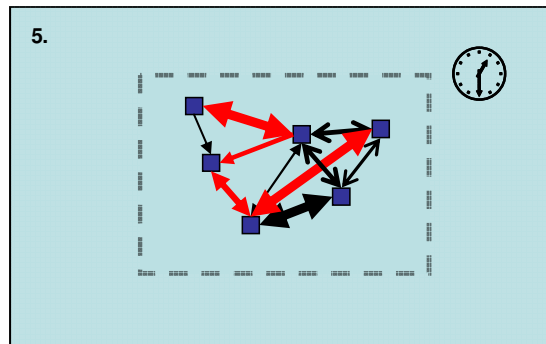
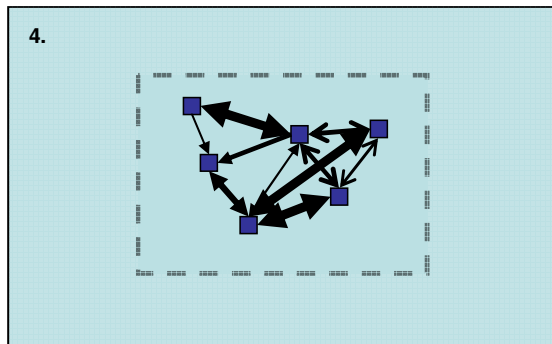
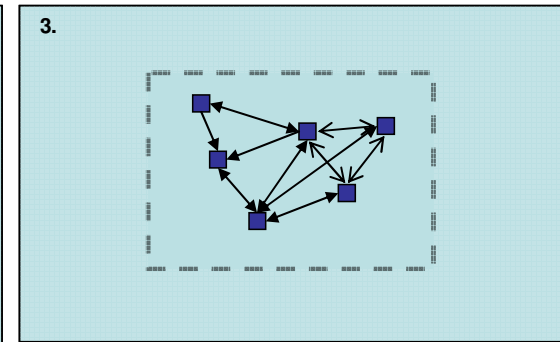
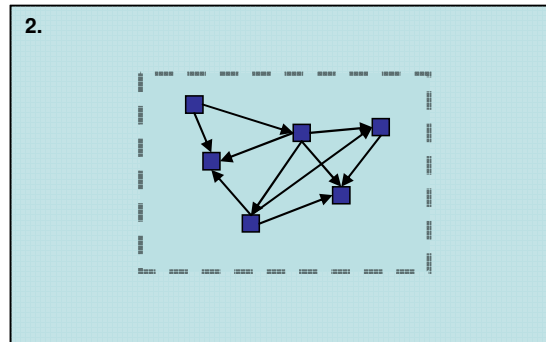
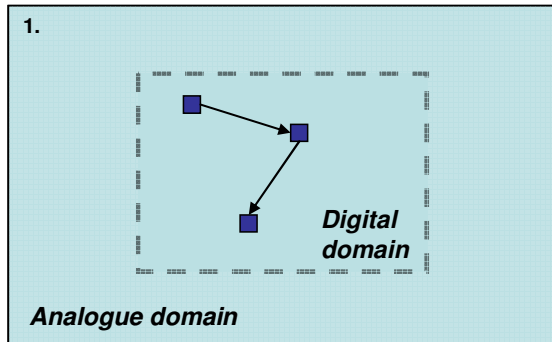
# Outline

- Background
- Case studies (C2C market places, social networking sites)
- IaD methods (user-centric, site-centric, network-centric)
- Lessons learned
- Added value IaD
- Implications for statistical agencies & policy makers
- Discussion

# Background

- Policy need to better understand phenomena associated with the Emerging Digital Economy (EDE)
- These phenomena are only partly captured in “established” statistics
- Notion of “digital footprints”
- With the advance of digitalization the scope of IaD-methods continuously expands

# Increasing scope of IaD methods



# Research questions

Two key questions:

- 1. Identify new data and indicators derived directly from the Internet and describe new phenomena associated with EDE**
- 2. To explore and assess the usefulness of the various IaD methods for deriving new, extra and substitute data for the EDE**

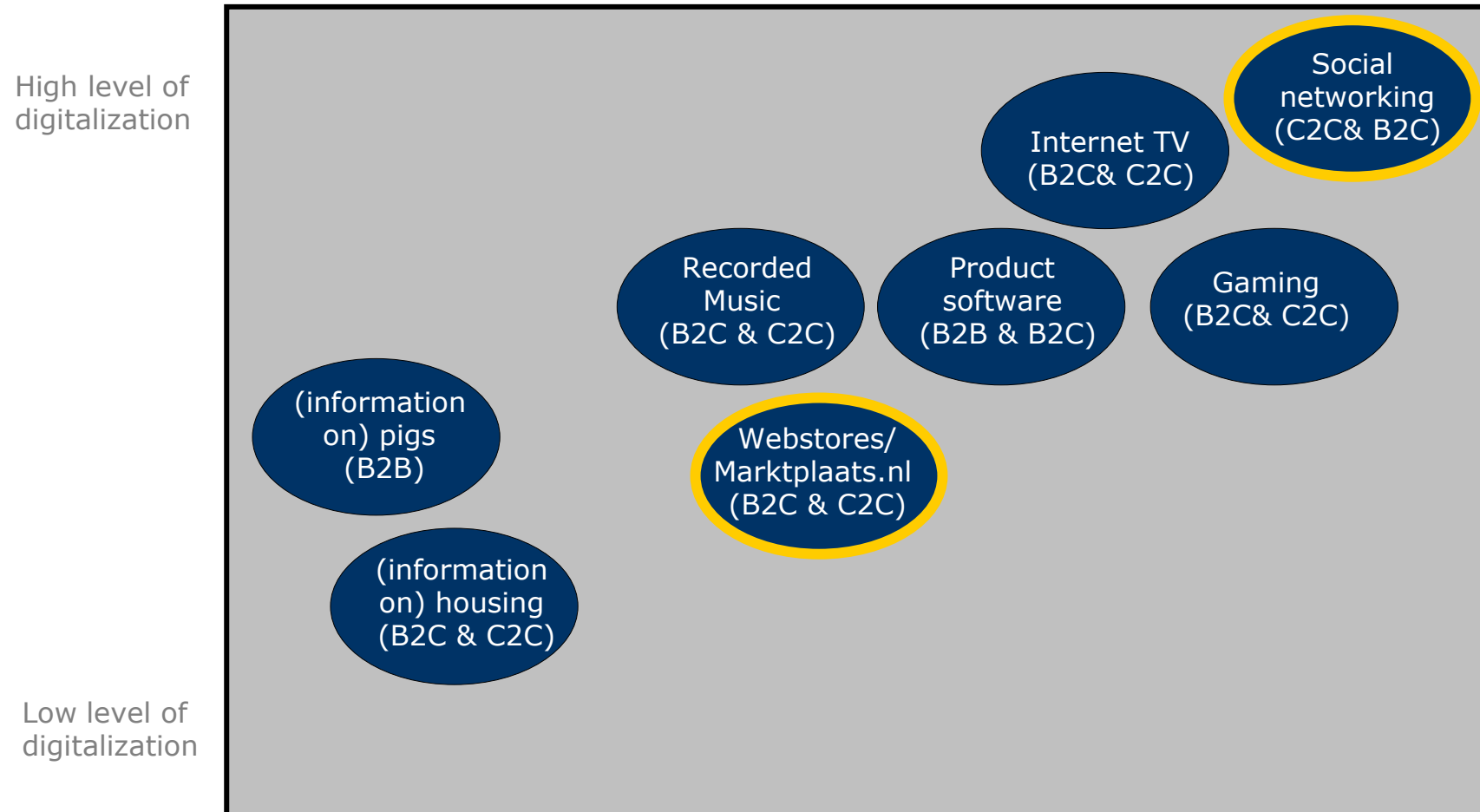
## Activities performed

- Analysis of international sources using IaD (desk research)
- Conceptualisation (**building blocks**)
- **8 Case studies** using fixed format, including some experiments
- **Typology IaD methods**: user-, network- and site-centric measurements
- Usability of spiders/web **crawlers/deep packet inspection**
- Overall analysis & reporting
- Contribution to CBS publication on the Digital Economy
- Presentation at OECD Working Party on Indicators for the Information Society (WPIIS), Paris, April 30 2008
- International consultation round (ongoing)

# Practical reach: covering both "old" & "new" economy

	<b>"Old economy" (established) markets &amp; phenomena</b>	<b>"New economy" (emerging) markets &amp; phenomena</b>
<b>Established (analogue, mostly invasive) data collection methods</b>	<b>(1)</b> ICT-investments in NACE-measured through a postal survey	<b>(2)</b> New media use by final users through a survey among a panel of households
<b>Internet-based (digital, mostly non-invasive) data collection methods</b>	<b>(3)</b> Prices of pigs traded over electronic market places measured through a site centric measurement	<b>(4)</b> Share of illegal content in P2P traffic as measured through a network centric measurement

# Selection of cases



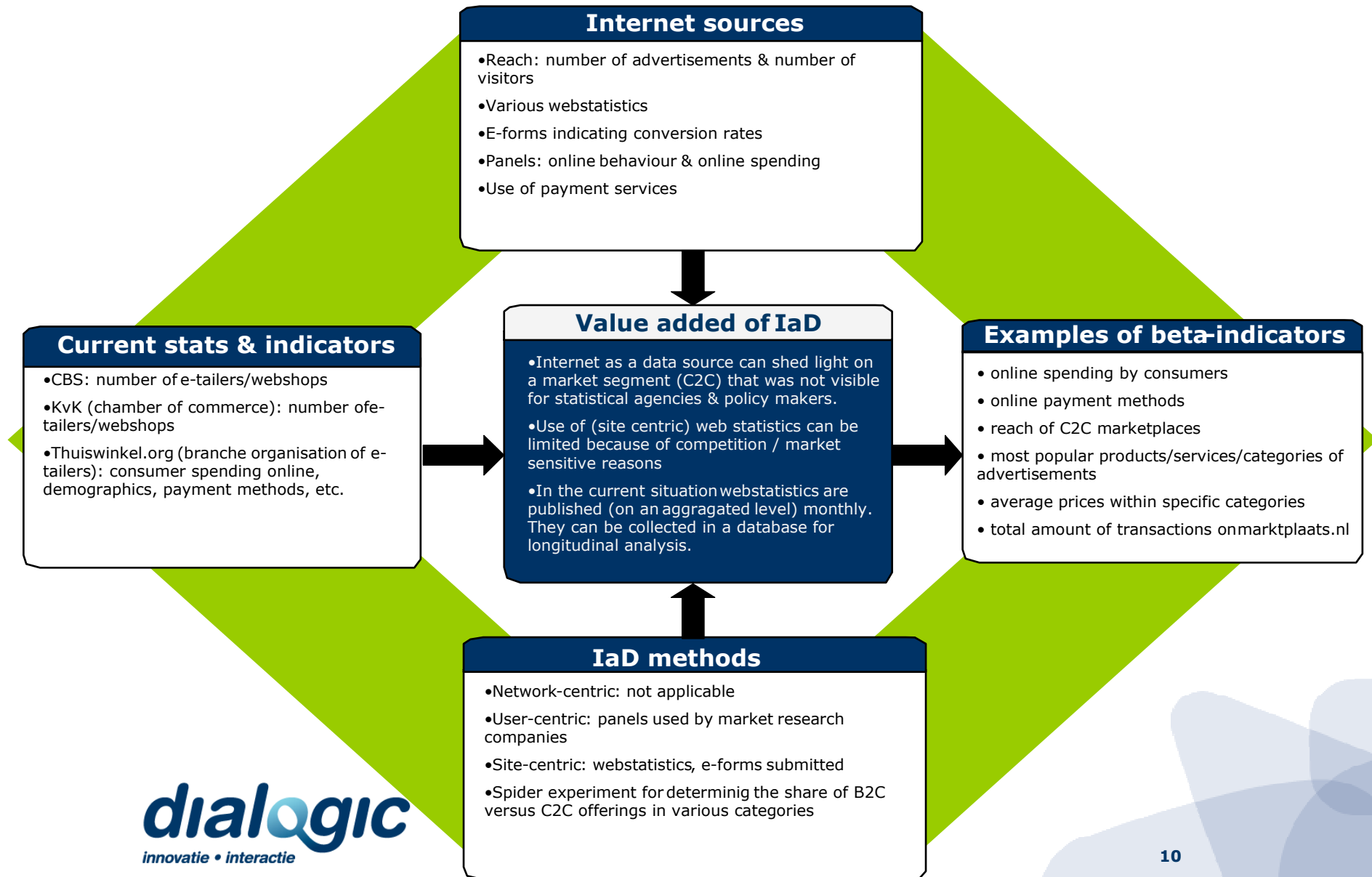


# Market places

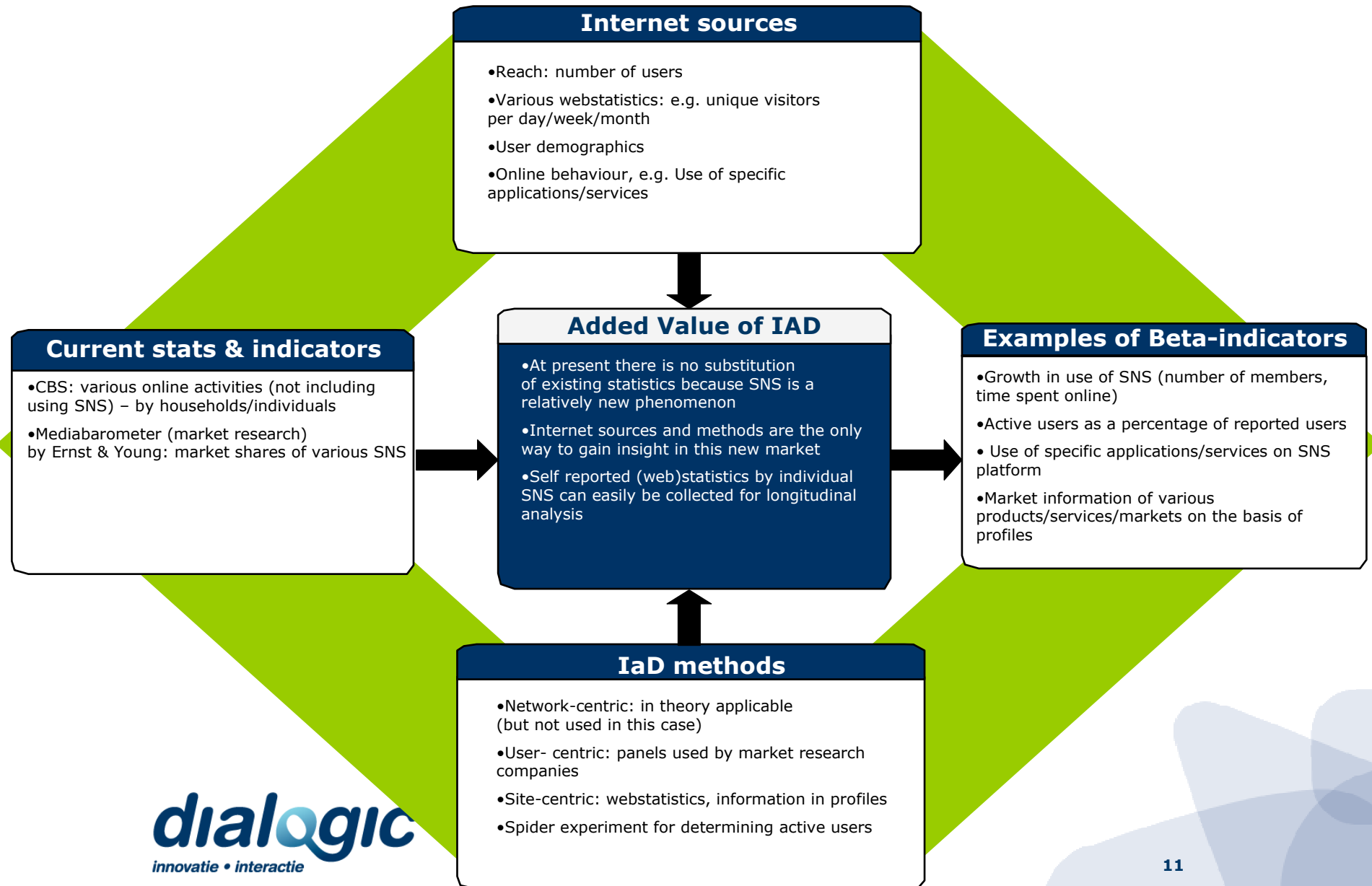
## Building block: Digital concentration points

	INFORMATION	ORDERING	PAYMENT	FULFILLMENT & LOGISTICS	
B2B	B2B online music marketplace	The 'Big Four' (Sony BMG, EMI, Universal, Warner)			
		Wholesaler of music on CD's and MP3's			
B2C	Physical music stores, e.g. Free Recordshop, Music Store, Media Markt				
	Online music stores, e.g. Bol, Proxis				
	Online sale of ring tones, e.g. Jamba, Boltblue				
	Online sale (legitimate) of MP3, e.g. iTunes				
			Provider of online financial services, e.g. Paypal, Ideal	Charts, e.g. Top-40, Top-50	
				Postal services (UPS, TNT)	
C2C	Online marketplaces, e.g. eBay, marktplaats, speurders				
	Internet search engines, e.g. Google, Yahoo				
	Torrent trackers, e.g. TorrentSpy, The Pirate Bay				Torrent trackers e.g. TorrentSpy, TPB
	Sociale networks, fan sites, LastFM,	Online storage, e.g. MegaUpload, RapidShare		Online storage e.g. MegaUpload, RapidShare	

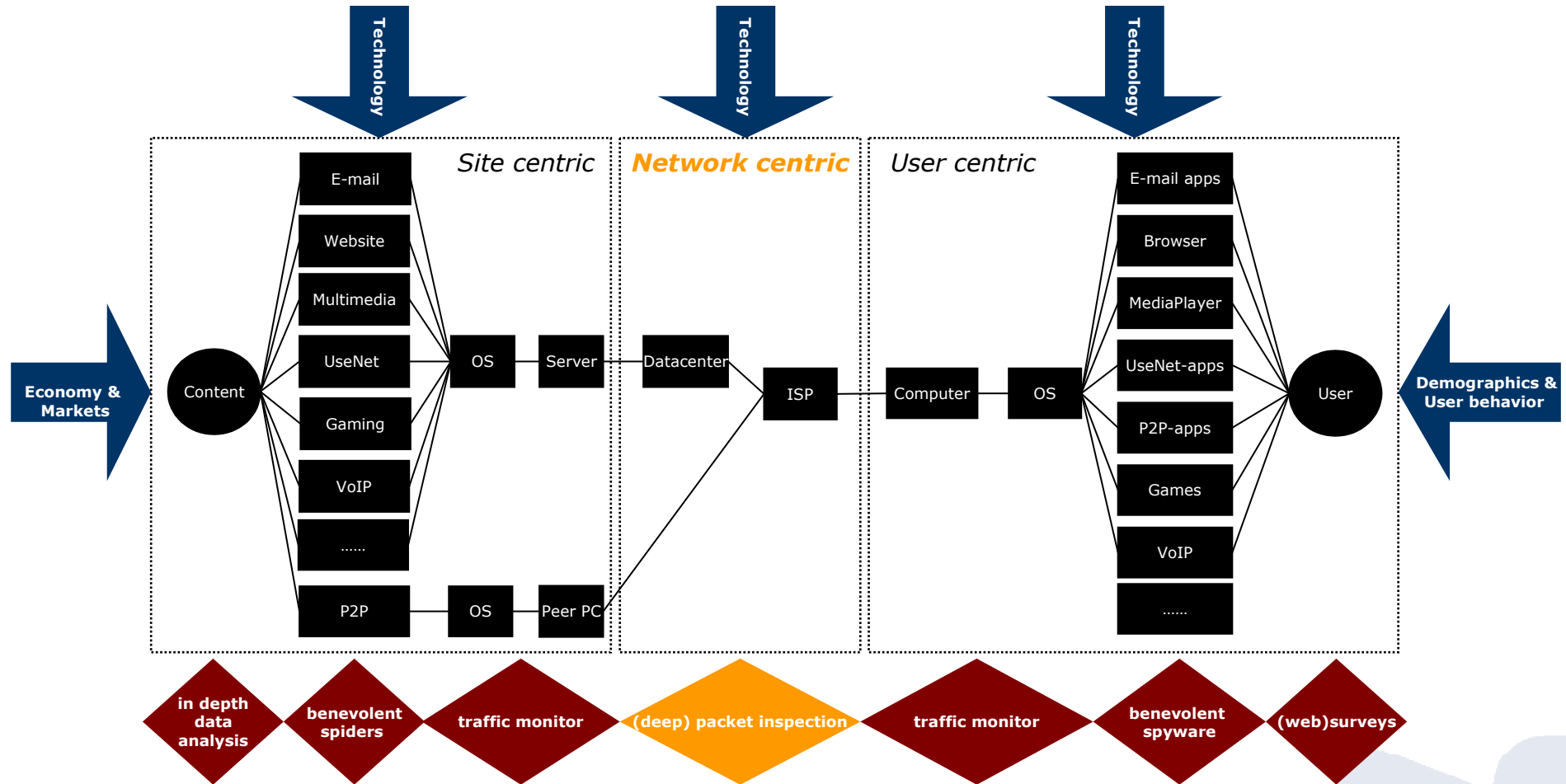
# Case (1): C2C online market places



# Cases (2): Social networking sites



# A typology of IaD-methods



Network-centric/Deep packet inspection

# Usability & disadvantages of IaD methods

Method	Disadvantages	When to use?
<b>Spiders</b>	<ul style="list-style-type: none"> <li>-Custom-made for every application</li> <li>-Feasible if the set of applications is limited and stable</li> <li>-Owner of an application can hinder being spidered</li> </ul>	<ul style="list-style-type: none"> <li>-When insight in content is needed</li> </ul>
<b>DPI@ISP</b>	<ul style="list-style-type: none"> <li>-The data does not allow generalisation</li> <li>-Difficult to obtain insight in content</li> <li>-Privacy of users can be threatened</li> <li>-Very hard to find ISPs willing to cooperate</li> </ul>	<ul style="list-style-type: none"> <li>-When a full scope of internet traffic is needed</li> <li>-When strong 'sympathy effects' are present</li> <li>-When very small effects and / or trends real time have to be identified</li> </ul>
<b>Traffic Monitor at OS</b>	<ul style="list-style-type: none"> <li>-Measurements are relatively shallow</li> <li>-A (costly) panel is needed</li> <li>-Illegal and shameful behaviour not measured correctly</li> <li>-The limited panel size makes is hard to find small effects</li> </ul>	<ul style="list-style-type: none"> <li>-When insight in user behaviour on all applications is needed.</li> </ul>
<b>Spyware</b>	<ul style="list-style-type: none"> <li>-Spyware needs to be custom-made</li> <li>-Measurements are relatively shallow</li> <li>-A (costly) panel is needed</li> <li>-Illegal and shameful behaviour not measured correctly</li> <li>-The limited panel size makes is hard to find small effects</li> </ul>	<ul style="list-style-type: none"> <li>-When insight in user behaviour on a certain application is needed.</li> </ul>

# Statistical pros and cons

IaD-method	Robustness (internal validity)	Representativity (external validity)	Transparency	Longitudinal use
<b>Spyware and Traffic Monitor</b>	<b>High</b> but underestimates illegal behaviour	<b>High.</b> depends on panel.	<b>Very High.</b> Like conventional surveys.	<b>High.</b> Sometimes changes in software
<b>DPI at ISP</b>	<b>High.</b> But advanced users can hinder DPI	<b>Low.</b> user characteristics are usually unknown.	<b>Very low.</b> Non-disclosure agreements	<b>Medium.</b> Small changes in infrastructure have major implications.
<b>Benevolent spiders</b>	<b>Low-medium.</b> Differences between 'websites' hinder measurement. structural bias  Underestimates illegal content	<b>Varies.</b> High in Concentrated markets. Low in fragmented market	<b>Medium.</b> OS Spiders	<b>Low.</b> Continuous changes 'websites'

## Lessons learned (1): overall lessons

1. IaD helps in signaling new trends, developments and phenomena
2. Mix of IaD methods applied will vary widely between markets and industries
3. The notion of digital footprints is not limited to digitalized products and services, but applies to a wider set of markets and industries





## Lessons learned (2): overall lessons

4. Information on goods and services represents an economic value in itself
5. Industry itself has already started to mine digital footprints
6. Markets associated with the emerging digital economy are more fuzzy and diffuse: traditional industry sector boundaries are getting problematic
7. In some markets user generated content (already) starts to mix with traditional economic production (online music, internet TV, housing) → barriers between the social and economic realm are blurring (SNS, C2C marketplaces)



## Lessons learned (3): overall lessons

8. Technical and practical availability of digital sources for third parties may differ considerably
9. Added value of using IaD may be higher in newly developing markets
10. IaD can shed light on the darker side or grey zone of the emerging digital economy



## Added value IaD (1)

- Relevant data source for policy-makers, researchers & statisticians, market research firms, industrialists and trade organizations
- Provide insight into markets and phenomena in areas where the statistical agencies have no established statistics available
- The potential of IaD methods for substitution of existing indicators and statistics should not be overestimated
- IaD-methods may lead to beta-indicators for the EDE → trade off between “early warning quality” indicators vs. lower statistical quality.

## Added value IaD (2)

Added value of using IaD is highest when:

- value chains are highly digitalized;
- online activities are the subject of research;
- subjects of research are highly dynamic and/or real time information about the subject is required.
- markets are dominated by a few players;
- market players are very transparent;
- markets are highly regulated (→ e.g. high quality registers)
- administrative tasks are labour intensive (scope for reducing administrative burden).

# Implications for statistical agencies

- Statistical agencies need to better capture phenomena associated with EDE – if they do not perform IaD methods themselves they should at least guarantee the quality of the statistical data/indicators that are generated by private firms (that are already filling the gaps that are left by public agencies)
- They are well positioned to play a key role in the switch to IaD-methods as they have:
  - scale and expertise for developing and collecting statistical indicators (sunk costs);
  - possibility to validate data and indicators derived from IaD measurements using regular statistics;
  - possibility to guarantee privacy if needed;
  - a judicial status they might want to use to enforce co-operation of data providers;
  - the international network for international benchmarking, exchange of expertise and setting standards and developing international guidelines.

# Implications for policy-makers

Various options to spur the further experimentation & use of IaD:

- new “beta statistics” publication on the emerging digital economy
- create a network of researchers, market research agencies, policy makers and statisticians
- establish a clearinghouse for Internet statistics.
- support statistical agencies to pro-actively experiment & use IaD-methods
- start exploratory talks with organisations and companies that can contribute to this R&D network
- governments themselves can anticipate on the use of digital sources for statistical purposes when developing or implementing their own registers and ICT projects

... to be taken up adopting an international perspective

# Discussion

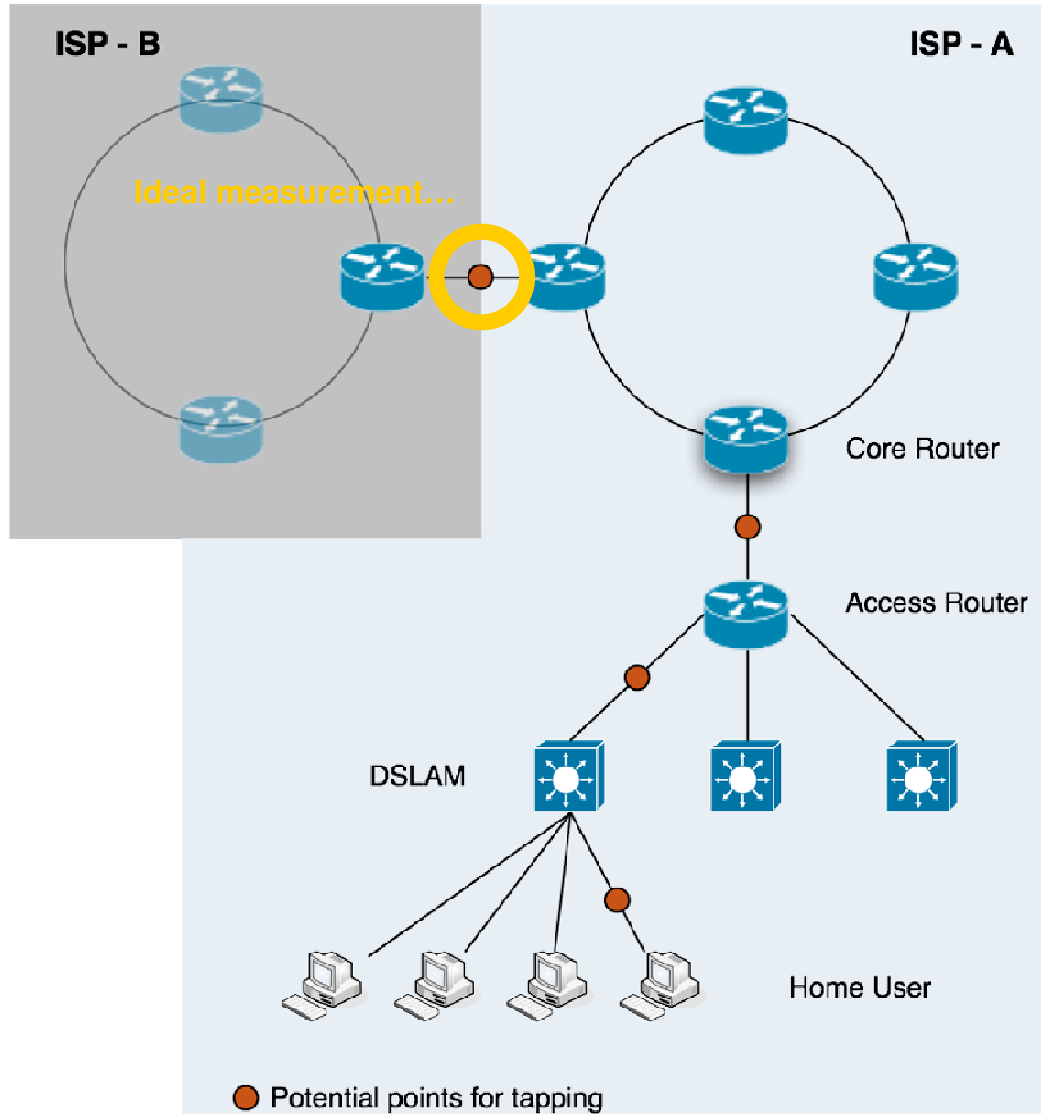
- Faustian bargain: trade-off between efficiency, objectivity, timeliness and cost-effectiveness on the one hand and validity and privacy on the other hand
- Sometimes there are simply no alternatives to the use of IaD methods
- How we see these beta-indicators:
  - There is a new category of beta-statistics that are especially suitable to pick up early trends in the EDE.
  - We should assess the practical and statistical quality of these statistics :
    - If these statistics do not meet the quality standards of beta-statistics they should be dropped.
    - The quality of the remaining beta-statistics should be improved (e.g., definitions, standardization, more sophisticated indicators).
    - Eventually, some beta-indicators could be promoted to the league of alpha-statistics.

<http://www.ipoque.com:80/website.html>

The port is part of a network address and normally hard coded into the application

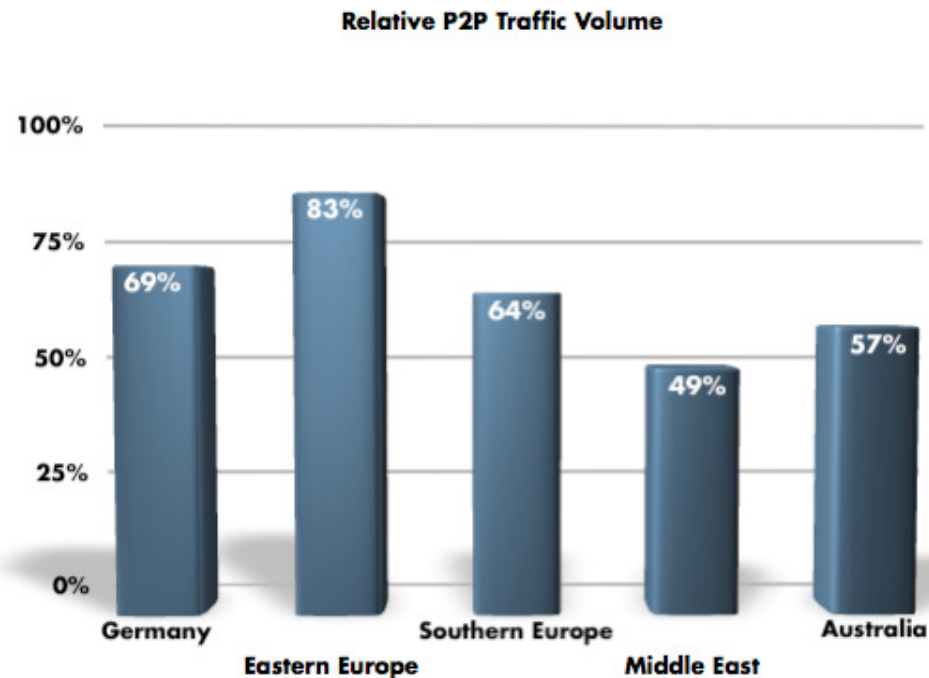
- Port based traffic classification:
  - does not work any longer
  - modern applications, like Skype or Instant Messengers are not bound to dedicated ports
- Deep Packet Inspection:
  - classification of network traffic based on unique application signatures





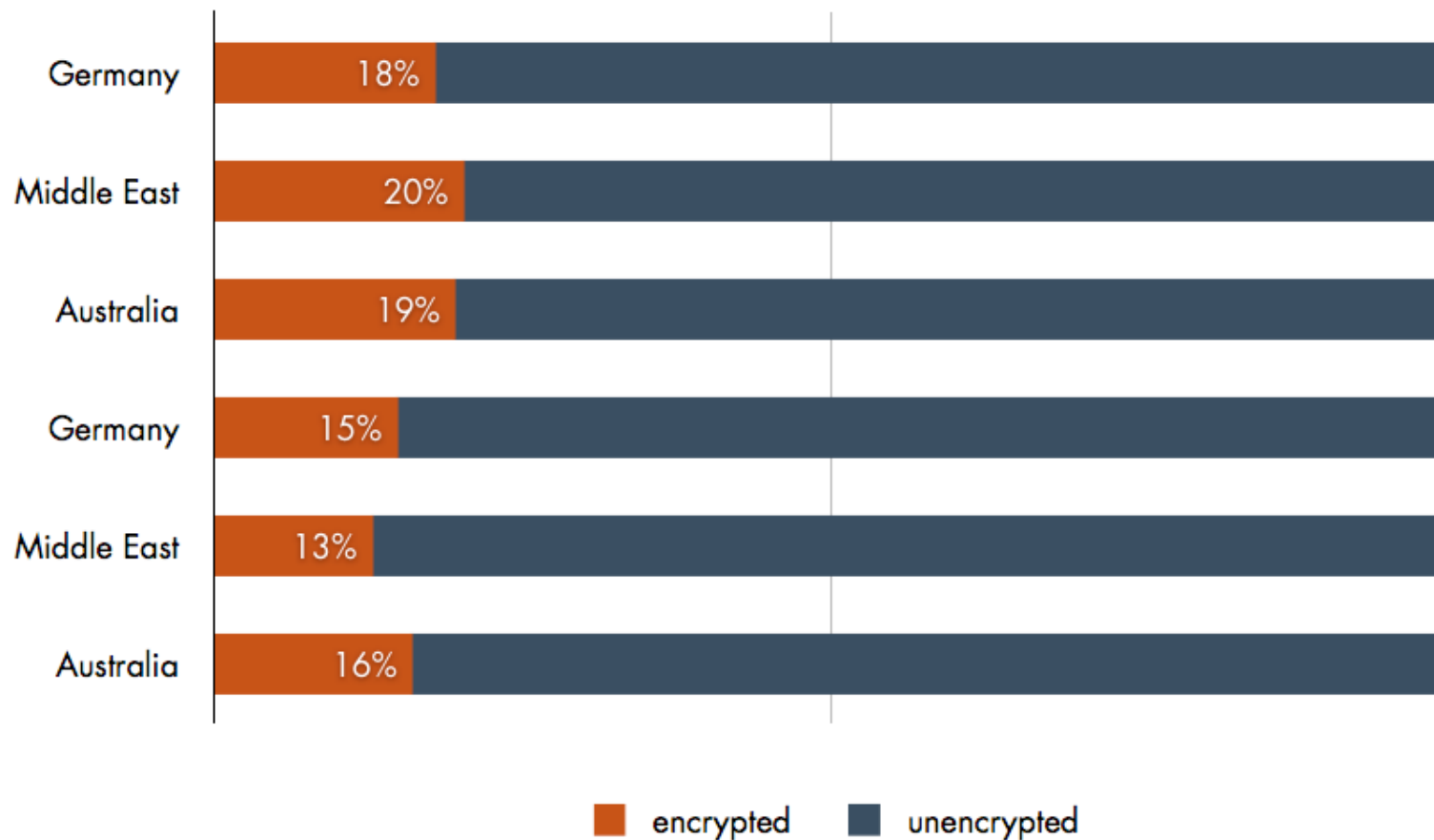
- Snapshot of the current state of the Internet
  - 18 monitoring sites at ISP (13) and Universities(5)
  - 5 regions
    - Southern Europe, Australia, Germany, Eastern Europe, Middle East
  - 3 Petabytes analyzed traffic
  - representing more than 1m people
  - data taken from the PRX Traffic Manager, installed at customers
  - not representative but a good estimation of
  - “What happens in the Internet”
  - Not just P2P, also VoIP, Skype, IM, Video Streaming, DDL

## More than 50% of the Internet traffic - worldwide

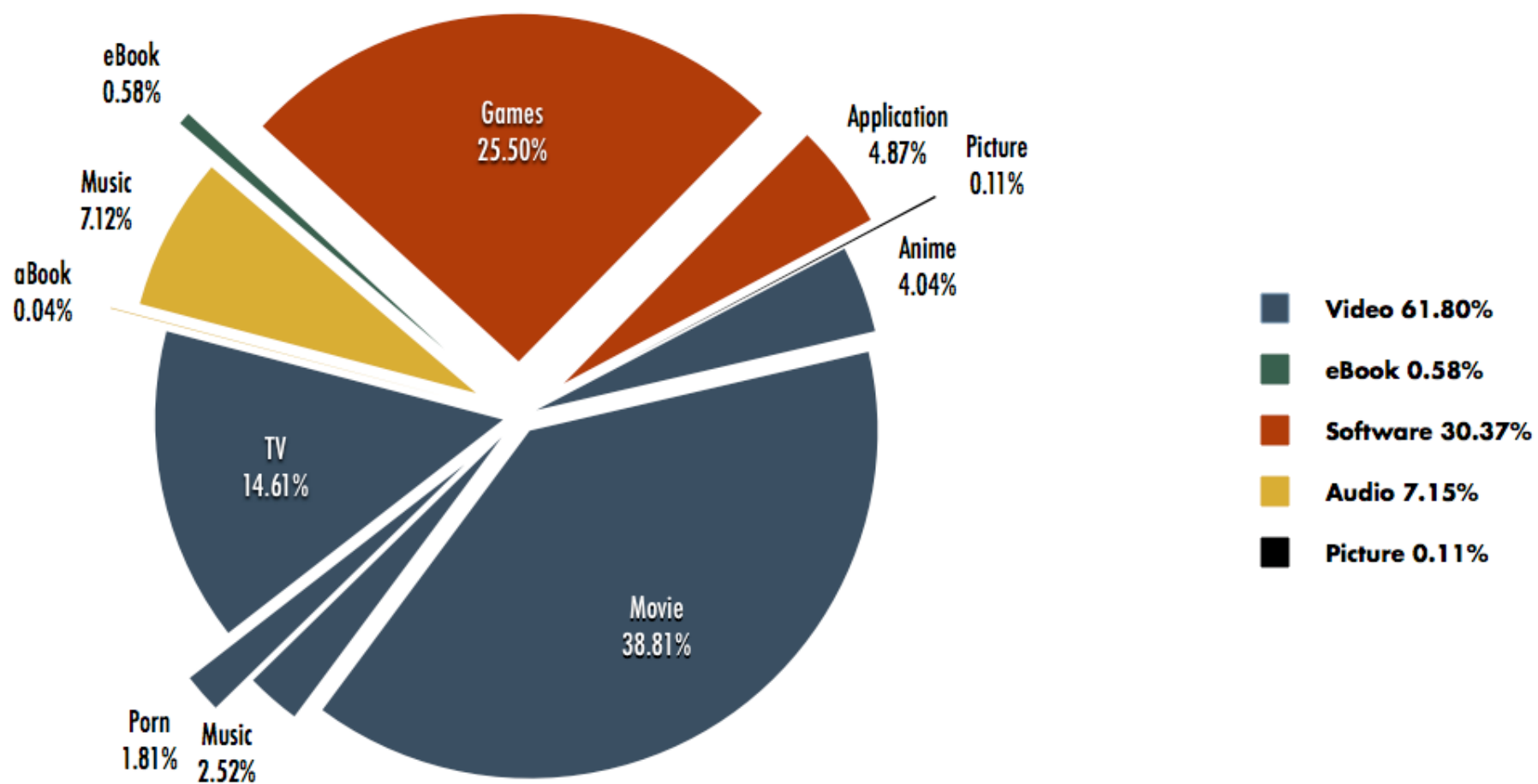


Source: ipoque Internet Study 2007

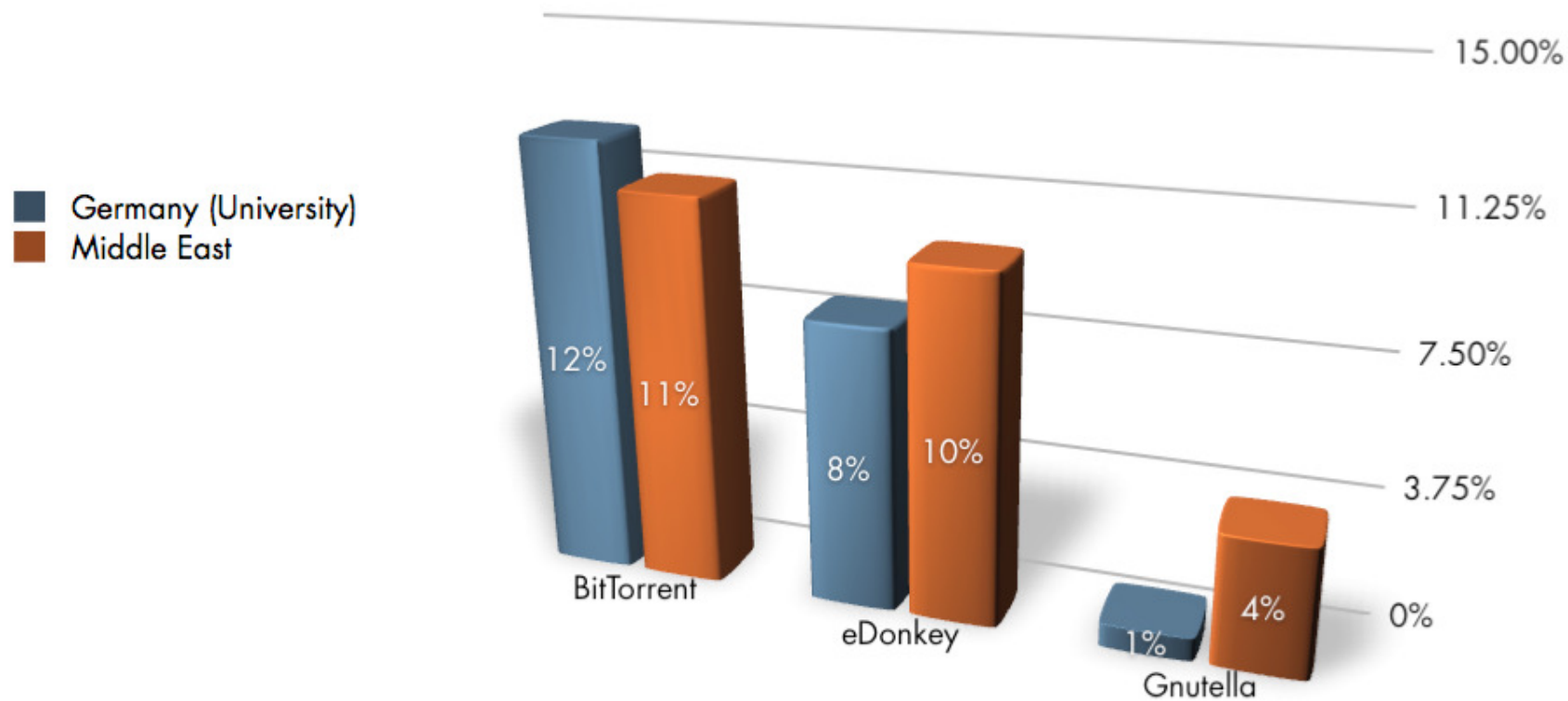
### Proportion of Encrypted P2P Traffic



**Traffic Volume per Content Type  
Southern Europe, BitTorrent**



Relative User Numbers per P2P Protocol



## Top 5 per Content Type (Southern Europe, BitTorrent)

### Video

1. Movie Next 2007
2. Movie The Simpsons Movie(Spanish)
3. Movie Shooter
4. Movie Evan Almighty
5. Movie Premonition

### Music

1. Music Bob Dylan-Blues-2006-MTD
2. Music Da Weasel 2007 Amor Escarnio e ....
3. Music Celine Dion 2007 D'elles
4. Music Bob Dylan - Live at the Gaslight 1962  
[2005]
5. Music Maroon 5 -It Won't be soon bevor long

### Software

1. Application K-Lite Mega Codec Pack 3.3.5
2. Games Football Manager 2007
3. Application Nero 7 Ultra Edition
4. Application Adobe Photoshop CS3
5. Games SilkRoad v1.110 Europe Legend 1

### eBooks

1. eBook Muscle & Fitness 101 Workouts
3. eBook Muay Thai - The Art of Fighting
4. eBook All Social Interactions Books
5. eBook Get the Dream Job- Cover letter Secrets ...
6. eBook tomtom map 6 75 ES and PT



# Q & A

Christiaan Holland  
+31(0)30 2150582  
holland@dialogic.nl  
www.dialogic.nl

Robbin te Velde  
+31(0)30 2150591  
tevelde@dialogic.nl  
www.dialogic.nl

Further information:

[www.ez.nl](http://www.ez.nl)

→ search on "iad english"