# Open data ecosystems – wishful thinking or successful business?

**PROFESSOR PER RUNESON**   ✉ **PER.RUNESON@CS.LTH.SE**   🐦 **@SOFTENGRESGRP**

# Data science

…is an **interdisciplinary** field that uses **scientific methods**, processes, algorithms and systems to **extract knowledge** and insights from noisy, **structured and unstructured data**, and apply knowledge and actionable insights from data across a broad range of **application domains**. [Wikipedia]

But where does the data come from?

Do you have enough quality data?

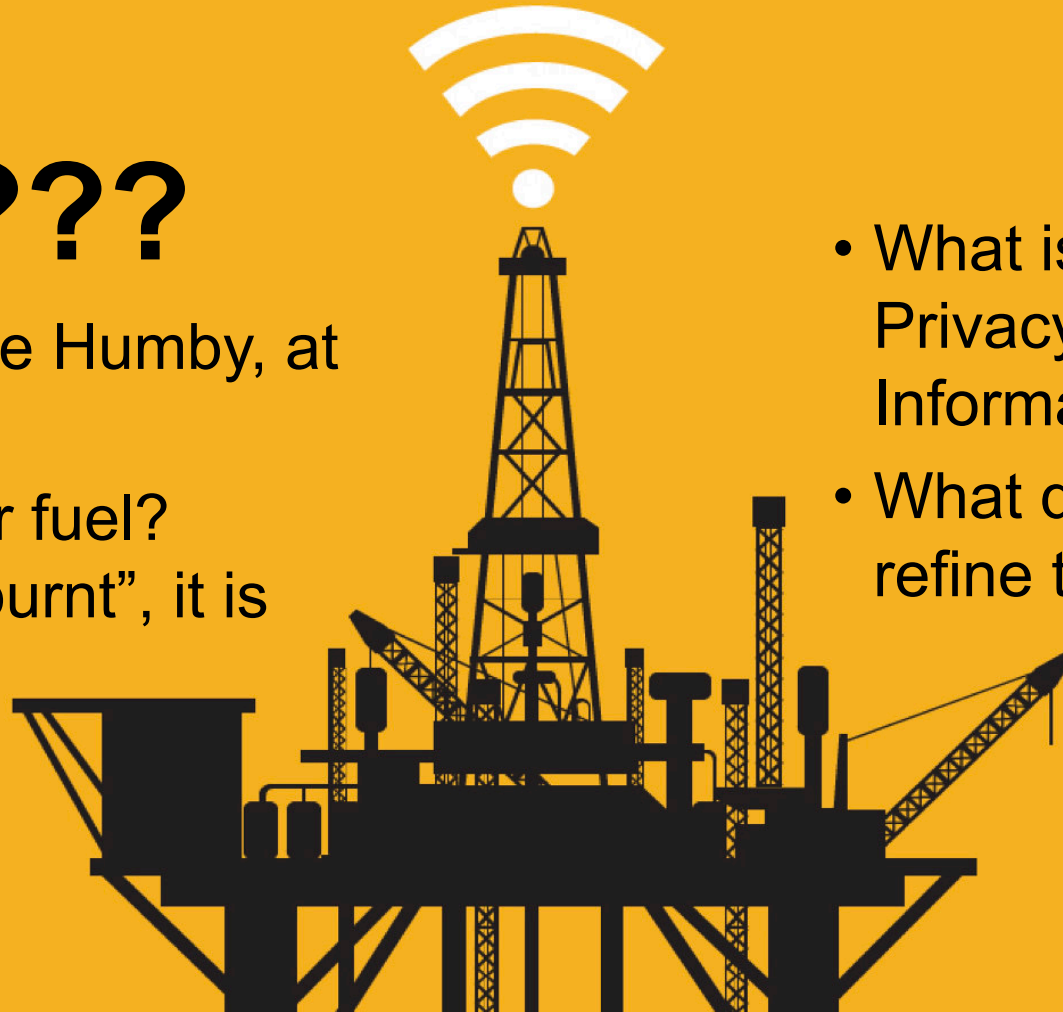Can you afford maintaining the data?

# Data is the new oil!

## Is it ????

- Claim by Clive Humby, at Tesco, 2006
- Lubrication or fuel? Data is not "burnt", it is non-rival

- What is data pollution? Privacy intrusion? Information leakage?
- What does it cost to refine the oil?

Gerd

# Example biomedicine

"For a typical biomedical data resource, the cost of simply keeping the data is only a small fraction of the total cost of data management. The remainder is largely the cost needed to support the finding, accessing, interoperating and reusing of the data — a cost that is widely under-appreciated."

# data for road safety

## Safety Related Traffic Information Ecosystem: Data for Road Safety
Live Vehicle, Crowd and Infrastructure Data improving road safety across Europe

Significantly improving road safety across Europe for all road users requires the mass involvement of vehicle manufacturers, traffic information service providers, automotive suppliers and public authorities. Such a level of participation will be necessary to ensure the pace and critical mass of safety data required for comprehensive safety related traffic information services.

**Update July 2021**

Privacy Statement – Data for Road Safety – 6 July 2021  ›

# Example maps

# Data challenges and opportunities

- Costs for *data maintenance*, *quality assurance* and *annotation* is a challenge

- Data will gradually become *commodity* for some functionality

Open data ecosystems?



Lundell *et al.* Commodification of Industrial Software: A Case for Open Source, *IEEE Software*, 26(04):77-83, 2009.  https://doi.org/10.1109/MS.2009.88

# What is unique?



Achieving Simplicity with the Three-Layer Product Model

Jan Bosch, *Chalmers University of Technology, Sweden*

| | | |
|---|---|---|
| **Innovation** | External sources | **Novelty focus** |
| ↓ Productize | | |
| **Differentiation** | | **Value focus** |
| ↓ Commoditize | | |
| **Commodity** | | **Cost focus** |

LUND UNIVERSITY

# Data sharing?

"Value comes from data being brought together, and that requires organizations to let others use the data they hold"

Regeringskansliet

https://www.regeringen.se/informationsmaterial/2021/10/data--en-underutnyttjad-resurs-for-sverige/

Sök på regeringen.se

**Sök**

Pressmeddelande från Infrastrukturdepartementet

# Ny nationell strategi ska göra Sverige ledande i delning av data

Publicerad 22 oktober 2021

**Genväg**

> Temasida: Data – en underutnyttjad resurs för Sverige

## Insatsområde 2: Öppen och kontrollerad datadelning

Mål 2023: Statliga myndigheter och statliga företag har en god förmåga att dela data både på ett öppet och kontrollerat sätt. Svenska företag har en god förmåga att dela data och är delaktiga i utvecklingen av och kan utnyttja de uppbyggda datamarknaderna. Offentliga data, inklusive forskningsdata, ska där så är lämpligt, vara så öppna som möjligt och så stängda som nödvändigt.

innovation.

# Background and motivation
# Open source software practices

# Background – Open Source Software

- 1960/70's – software into the bargain
- 1980's – political movement
- 1990's – commercial (Linux)
- 2000's – databases (MySQL), Android
- 2010's – everywhere

https://dx.doi.org/10.1109/MC.2020.3041887



OPEN SOURCE EXPANDED

EDITOR **DIRK RIEHLE**
Friedrich Alexander-University of Erlangen Nürnberg;
dirk.riehle@fau.de

A Brief History of Free, Open Source Software and Its Communities

**Jesus M. Gonzalez-Barahona,** Universidad Rey Juan Carlos

Free, open source software (FOSS) has a long history, beginning with the origins of software itself, when the terms free software and open source software were not yet defined. Learning about the milestones of this history may help to understand FOSS today,

computers, with IBM being, by a large margin, the market leader. For all of them, software was just a companion to hardware: as long as you paid for maintenance, you had access to the software catalog of the manufacturer. User groups, such as SHARE (IBM) and the DECUS [Digital Equipment Corp. (DEC)] favored software sharing. To some ...
to 1970. ...

# Open source in mobile devices – 2011



**Fig. 1.** Worldwide smart-phone Market shares (%) by platform in 2009/2010 (Gartner, 2011)

## Open-Source Software Implications in the Competiti[ve] Mobile Platforms Market

Salman Qayyum Mian[1], Jose Teixeira[2], and Eija Koskivaara[3]

[1] Nokia Siemens Networks (NSN), Linnoitustie 6, 02600 Espoo, Finland
Salman.Mian@uta.fi
[2] Turku Center for Computer Science (TUCS), Joukahaisenkatu 3-5 B, 20520 Turku, Finland
Jose.Teixeira@tse.fi
[3] Turku School of Economics (TSE), Rehtorinpellonkatu 3, 20500 Turku, Finland
eija.koskivaara@utu.fi

**Abstract.** The era of the PC platform left a legacy of competitive strategies for the future technologies to follow. However, this notion became more complicated, once the future grew out to be a present with huge bundle of innovative technologies, Internet capabilities, communication possibilities, and ease in life. A major step of moving from a product phone to a smart phone, eventually to a mobile device has created a new industry with humongous potential for further developments. The current mobile platform market is witnessing a platforms-war with big players such as Apple, Google, N[okia,] Microsoft in a major role. An im[...]

# Triggers of Openness – why engage?

- Access to skilled workforce
- Faster development speed
- Low license costs and switching costs
- Flexibility in tool usage and adaptations
- Shared cost with the ecosystem
- Governing ecosystem

Feature Article

## How Companies Use OSS Tools Ecosystems for Open Innovation

**Hussan Munir**
Lund University

**Per Runeson**
Lund University

**Krzysztof Wnuk**
Blekinge Institute of Technology

*Abstract*—Moving toward the open innovation (OI) model requires multifaceted transformations within companies. It often involves giving away the tools for prod development or sharing future product directions with open tools ecosystems. Mo from the traditional closed innovation model toward an OI model for software development tools shows the potential to increase software development compete and efficiency of organizations. We report a case study in software-intensive compa developing embedded devices (e.g., smartphones) followed by a survey in OSS communities such as Gerrit, Git, and Jenkins. The studied branch focuses on develo Android phones. This paper presents contribution strategies and triggers for openn These strategies include avoid forking OSS tools, empower de the ecosystem, steer ecosystems thr differentiation

# Strategies for open tools

A theory of openness for software engineering tools in software organizations

Hussan Munir[*,a], Per Runeson[a], Krzysztof Wnuk[b]

[a] Department of Computer Science, Lund University, P.O. Box 118, SE-221 00 Lund, Sweden
[b] Software Engineering Research Lab, Blekinge Institute of Technology, SE-371 79 Karlskrona, Sweden

**Strategy**

|  | Cost saving | Inspirational |
|---|---|---|
| **Proactive** | **Lucrativeness** (Think tank) | **Leaders** (Growth through ecosystems) |
| **Reactive** | **Laggards** (Business as usual) | **Leverage** (Resource optimization) |

**Why**

LUND UNIVERSITY

Open source for data?
Data Ecosystems!!!

# A Data Ecosystem is...

- a network **community** with a common interest
- supported by a **technological platform**
- to **process data**
  - e.g., find, archive, publish, consume, or reuse
- collaboration on **the data and resources**
  - e.g., software and standards

Vi är Trafiklab

# Data Ecosystem Dynamics

# Data Ecosystem Roles



Level of influence
- High
- Low

Data consumers

Data Producers
- Platform provider
- Keystone members
- Passive members
- End-users

General contributors

# How open is open?

# Essential concepts of Open Data Ecosystems

## PER RUNESON, THOMAS OLSSON, JOHAN LINÅKER

## Open Data Ecosystems — An empirical investigation into an emerging industry collaboration concept

Per Runeson [a,*], Thomas Olsson [b], Johan Linåker [a]

[a] Department of Computer Science, Lund University, Lund, Sweden
[b] Systems Engineering, RISE Research Institutes of Sweden AB, Lund, Sweden

A B S T R A C T

... increasingly depending on data, particularly with the rising use of machine ... sources of data. Open Data Ecosystems (ODE) is an ... systems, similar to Open Source ... share

# Emerging data ecosystems

**JobTech**

- Labor market
- Job ads
- Public-driven
- Organization-centric

**ESS-CSDL**

- Industry 4.0
- Alarm data
- Business-driven
- Organization-centric

**RoDL**

- Automotive
- Traffic video
- Business-driven
- Consortium-based

external/
internal
pecuniary/
non-pecuniary

Ches-
brough

[business models]
[tool support]

[business driven]

[co-opetition]

F7. Maturity

F1. Value of data

[knowledge]

organization-centric
consortium-based
community-based

Dal
Bianco

F5. Competition

**[G4. Evolution]**

**[G1. Value]** —— F2. Value of collaboration

ODE

coreness
currentness
granularity
+degree of processing

Enders

[platform
ownership]

**[G3. Governance]**

F0. Type

platform provider
keystone members
passive members
end users

Naka-
koji

F4. Relationship

**[G2. Intrinsics]** — F6. Quality —— [standardization]

F3. Acquisition

F8. Legal

[transparency]

[meta-data]
[domain model]

[degree of
openess]

closed
shared
open

Coyle

[public-driven]
[business-driven]
[community-driven]

[licenses]
[privacy]

[liability]

Open Data Ecosystems – an empirical investigation into an emerging industry collaboration concept

LUND
UNIVERSITY

# Value

The value of data (F1) and the value of collaboration around the data (F2) are two sides of the same coin. One or the other may be the primary value, but they are highly intertwined.

LUND
UNIVERSITY

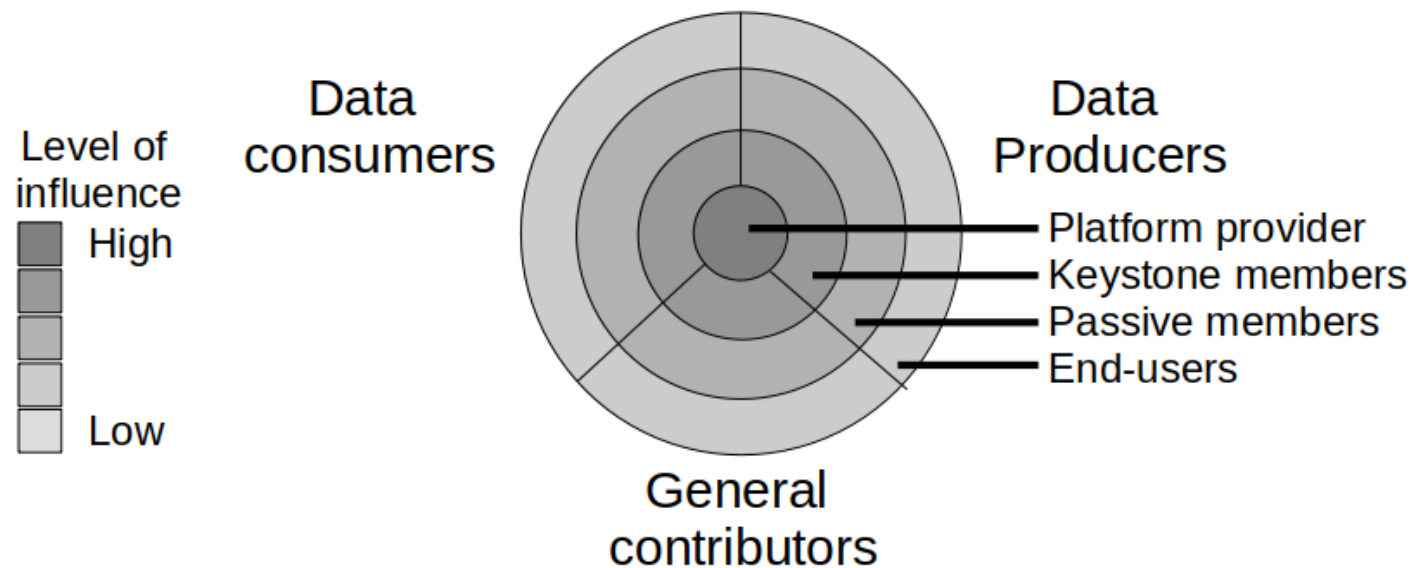# Intrinsics

Intrinsics,
or internal characteristics of data

- data type (F0)
    - coreness
    - currentness
    - granularity
    - degree of processing
- data quality (F6)
    - correctness
    - provenance
    - meta-data

- legal aspects (F8) is tightly connected to data, although they also connect to governance of the ODE.
    - licenses
    - privacy
    - liability

# Governance



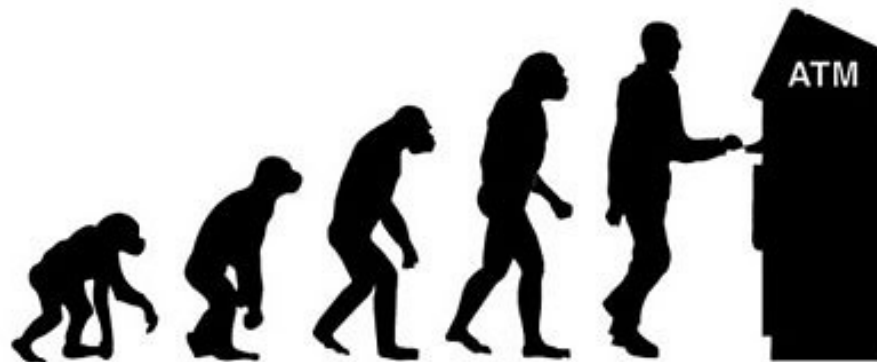There is a need for an independent platform provider to ensure trust

Initiation may be public-driven, business-driven, or community-driven

# Evolution

The concept of and strategies for open data ecosystems are still in their infancy

Need for knowledge:

– how to integrate ODEs into an organization's business model
– tools to support ODEs and enable data sharing should be developed and standardized

LUND
UNIVERSITY

# Findings for data ecosystems

**Value**

Focus on business value in the data or collaboration

**Intrinsics**

Data coreness, currentness and granularity
Standardize format and legal framework

**Governance**

Level of openness and platform ownership
Relationship and competition must co-exist
Data acquisition incentives

**Evolution**

Advance business models and tool support

# Recommendations for public platform providers

JOHAN LINÅKER, PER RUNESON

## How to Enable Collaboration in Open Government Data Ecosystems: A Public Platform Provider's Perspective

Johan Linåker and Per Runeson

Dept. of Computer Science, Lund University, Lund (SWE), johan.linaker@cs.lth.se, per.runeson@cs.lth.se

Abstract: Open Government Data (OGD) is an important driver for open innovation among public ... search highlights a need for improved feedback loops, collaboration, ... this study, we explore how public platform ... both in terms

# Open Governmental Data

Purpose: 1) Governance transparency, 2) Business innovation



Fig. 1. Elements of an open government data ecosystem derived from the literature.

## Innovation with open data: Essential elements of open data ecosystems

Anneke Zuiderwijk*, Marijn Janssen and Chris Davis
Faculty of Technology, Policy and Management, Delft University of Technology,

Open data ecosystems are expected to bring many advantages, such as stimulating ... an open data ecosystem ...
... attention has been given to what constitutes an open data ecosystems for enabling easy publ...
... tial elements of open data ecosystems ... scenario about the publicati...
... has been reviewed and ... tools and portals are available which togeth...
... utilized by open data providers and u...
... leasing and publishing o...
... ing, analyzing, ...
... to the ...

# Open Government Data ecosystems

**JobTech**

- Labor market
- Job ads
- Public-driven
- Organization-centric

**Trafiklab**

- Public transport
- Schedule, traffic
- Public-driven
- Consortium-based

**HSL DevCom**

- Public transport
- Schedule, traffic
- Public-driven
- Organization-centric

**City of Chicago**

- City governance
- All kinds of city
- Public-driven
- Organization-centric

# Recommendations for the public platform provider

Build community and trust

Maintain a clear vision

Be active and multi-functional

Build open communication

Develop open source software

Share data, other than your own

Adopt and promote open standards

# Open data ecosystems – wishful thinking or successful business?

Maybe data isn't the new oil?

It might be the new, renewable bio-energy
but we have to make it together

# More to read

The Journal of Systems & Software 182 (2021) 111088

Contents lists available at ScienceDirect

## The Journal of Systems & Software

journal homepage: www.elsevier.com/locate/jss

ELSEVIER

## Open Data Ecosystems — An empirical investigation into an emerging industry collaboration concept

Per Runeson [a,*], Thomas Olsson [b], Johan Linåker [a]

[a] Department of Computer Science, Lund University, Lund, Sweden
[b] Systems Engineering, RISE Research Institutes of Sweden AB, Lund, Sweden

ARTICLE INFO

1. Introduction

---

JeDEM

## How to Enable Collaboration in Open Government Data Ecosystems: A Public Platform Provider's Perspective

Johan Linåker and Per Runeson

Science, Lund University, Lund (SWE), johan.linaker@cs.lth.se, per.runeson@cs.lth.se

innovation among public

---

FOCUS: **GUEST EDITORS' INTRODUCTION**

# Collaborative Aspects of Open Data in Software Engineering

Johan Linåker, RISE Research Institutes of Sweden

Per Runeson, Lund University

Anneke Zuiderwijk, Delft University of Technology

Amanda Brock, OpenUK

Upcoming IEEE Software
January/February 2022

# More to come:
# B2B Data Sharing for Industry 4.0 Machine Learning

Prof. **Per Runeson**, PhD Student **Konstantin Malysh**, Software Engineering, LU
Prof. **Christian Kowalkowski**, PhD Student **Tanvir Ahmed**, Industrial Marketing, LiU

Kowalkowski

Ahmed

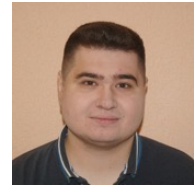## Business models (LiU)

Two disruptive and interrelated transformations:
1) **digitalization** changes sociotechnical systems,
2) **servitization** entails the shift from selling products to 'product-as-a-service' business models

## Collaboration tools (LU)

Git, Jenkins and Gerrit, provide a low-threshold entry o open source software (OSS). Data ecosystems need "an underpinning technological platform".

Runeson

Malysh

# WE ARE OPEN

per.runeson@cs.lth.se
http://www.lth.se/digitalth