# What is a text within the Digital Humanities, or some of them at least?

*Manfred Thaller, Universität zu Köln*

*Digital Humanities 2012, July 20th 2012*

# Information I

# Shannon

Claude Shannon: "A Mathematical Theory of Communication", Bell System Technical Journal, 1948.
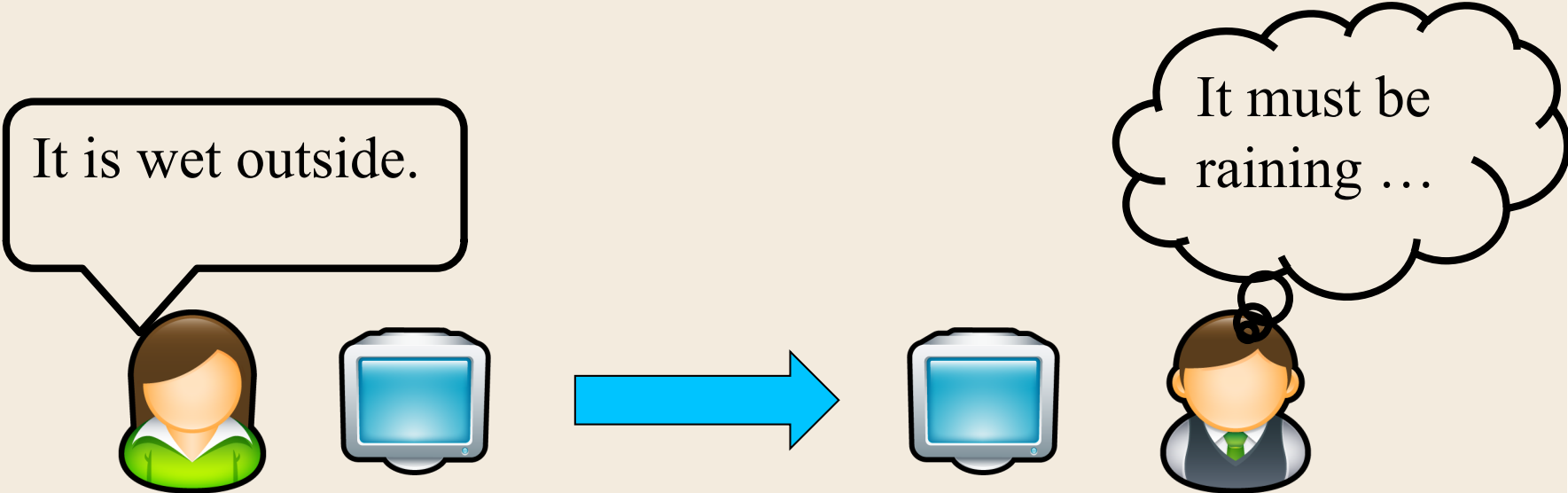
# Shannon

The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point.

*(Shannon, 1948, 379)*
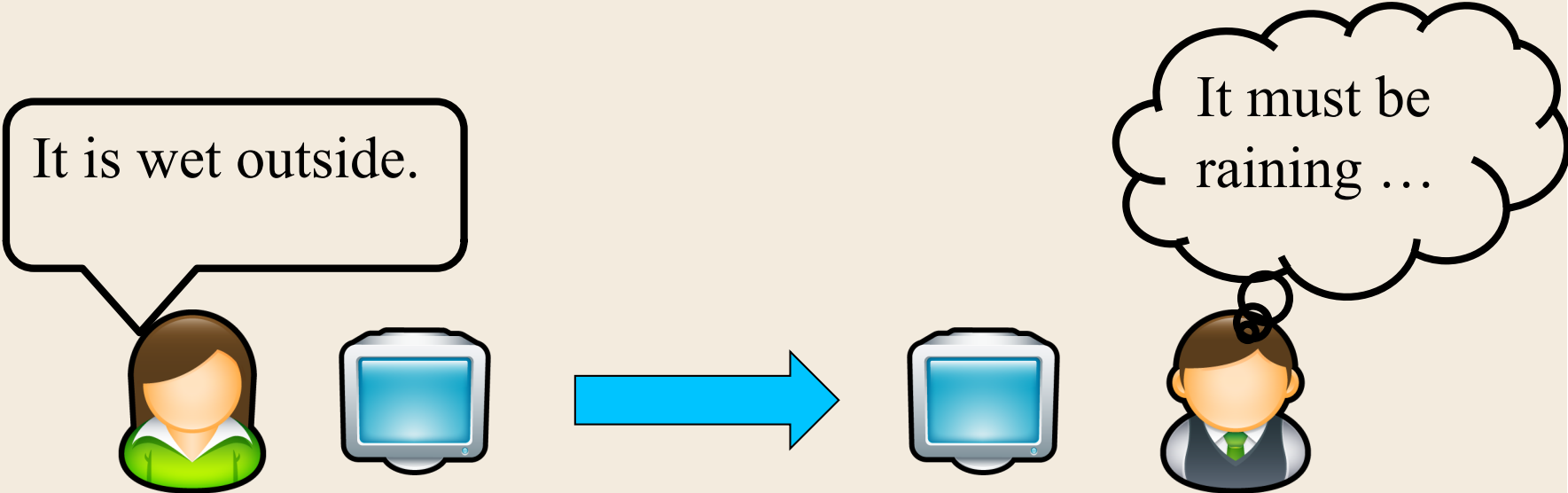
# Shannon

The fundamental problem of communication is that of reproducing at one point either exactly or approximately a message selected at another point. Frequently the messages have meaning; that is they refer to or are correlated according to some system with certain physical or conceptual entities. These semantic aspects of communication are irrelevant to the engineering problem.
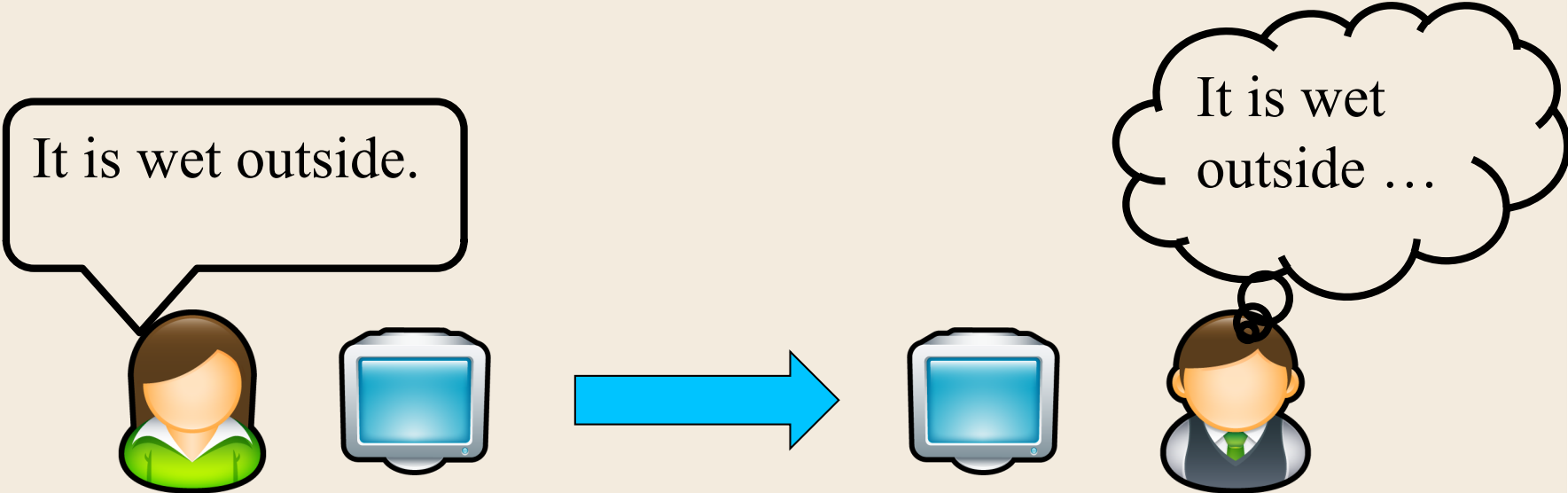
*(Shannon, 1948, 379)*

# Shannon

It is wet outside.

It must be raining …

# Shannon

# Shannon

# „Ladder of Knowledge"

# Information

# Data

# Data ➜ Information

*Data* are stored. E.g.:  *22°C.*

*Information* are data interpreted within a context:
"*In this lecture hall the temperature is  22°C*".

This context is fixed and identical for all recipients of information.

# Information ➔ Knowledge

*Knowledge* is the result of a more complex process.

E.g. the decision, derived from the room temperature of 22 ° centigrade, to get out of your jacket; or not.

This context is different between recipients of information.

# So …

| Data | Information |
|------|-------------|
| 22 ° C | 22 ° C in lecture hall M |
| 22 | 22 ° |
| '00000000 00010110' | 22 [ NOT ASCII { 0, 22 } ] |

# Langefors

Langefors "Infological Equation": original

# I = i (D, S, t)

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

# Information II

# Langefors

Langefors "Infological Equation": original

## I = i (D, S, t)

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

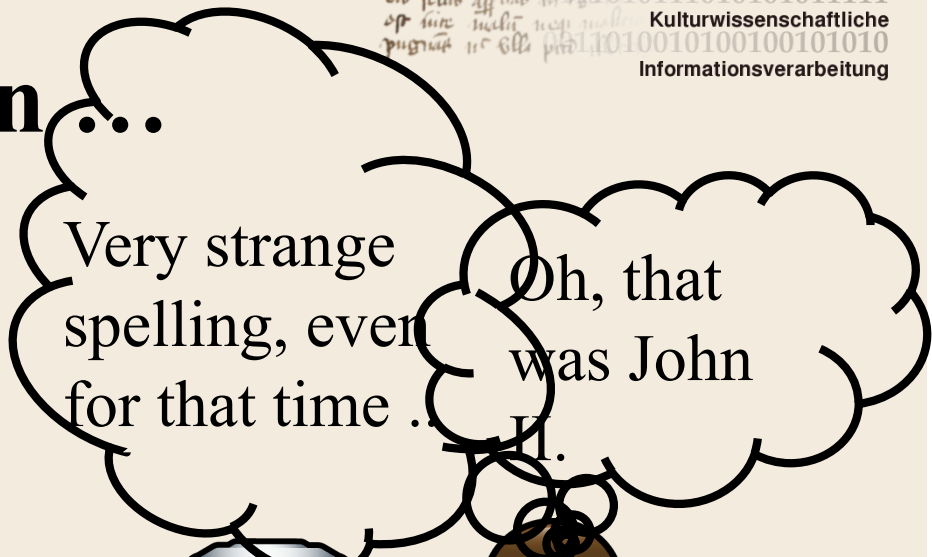Börje  Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

# Langefors

Langefors "Infological Equation": generalization 1

$$I_2 = i (I_1, S_2, t)$$

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

# Langefors

Langefors "Infological Equation": generalization 2

$$I_x = i (I_{x-1}, S_x, t)$$

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

# Langefors

Langefors "Infological Equation": generalization 3

$$S_x = s (I_{x-1}, t)$$

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge, s = knowledge generating process*

t ::= time

# Langefors

Langefors "Infological Equation": generalization 4

$$I_x = i (I_{x-\alpha}, S_{x-\beta}, t)$$

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

# Langefors

Langefors "Infological Equation": generalization 5

$$I_x = i\,(I_{x-\alpha},\, s(I_{x-\beta},\, t),\, t)$$

Börje Langefors, *Essays on Infology,* Studentliteratur: Lund, 1995

I ::= *Information*

i() ::= *interpretative process*

D ::= *Data*

S ::= *Previous knowledge*

t ::= time

# Remember …

| Data | Information |
|------|-------------|
| 22 ° C | 22 ° C in lecture hall M |
| 22 | 22 ° |
| '00000000 00010110' | 22 [ NOT ASCII { 0, 22 } ] |

# Changeable datatypes

int myVariable;

char myVariable;

temperature myVariable;


obj myVariable;

myVariable.useAsInt();

myVariable.useAsChar();

myVariable.addInterpretation(temperature,Centigrade);

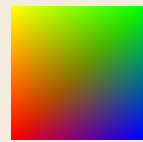# Langefors

Notes:

(1) If this is so, the assumption of Comp. Sci., that information is represented by structures on which algorithms operate, can be replaced by a more general understanding, according to which information is a state of a set of perpetually active algorithms.
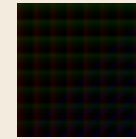
(2) Has that any practical meaning?
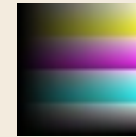
# A practical interlude

# Planets: the problem



► Photoshop ►
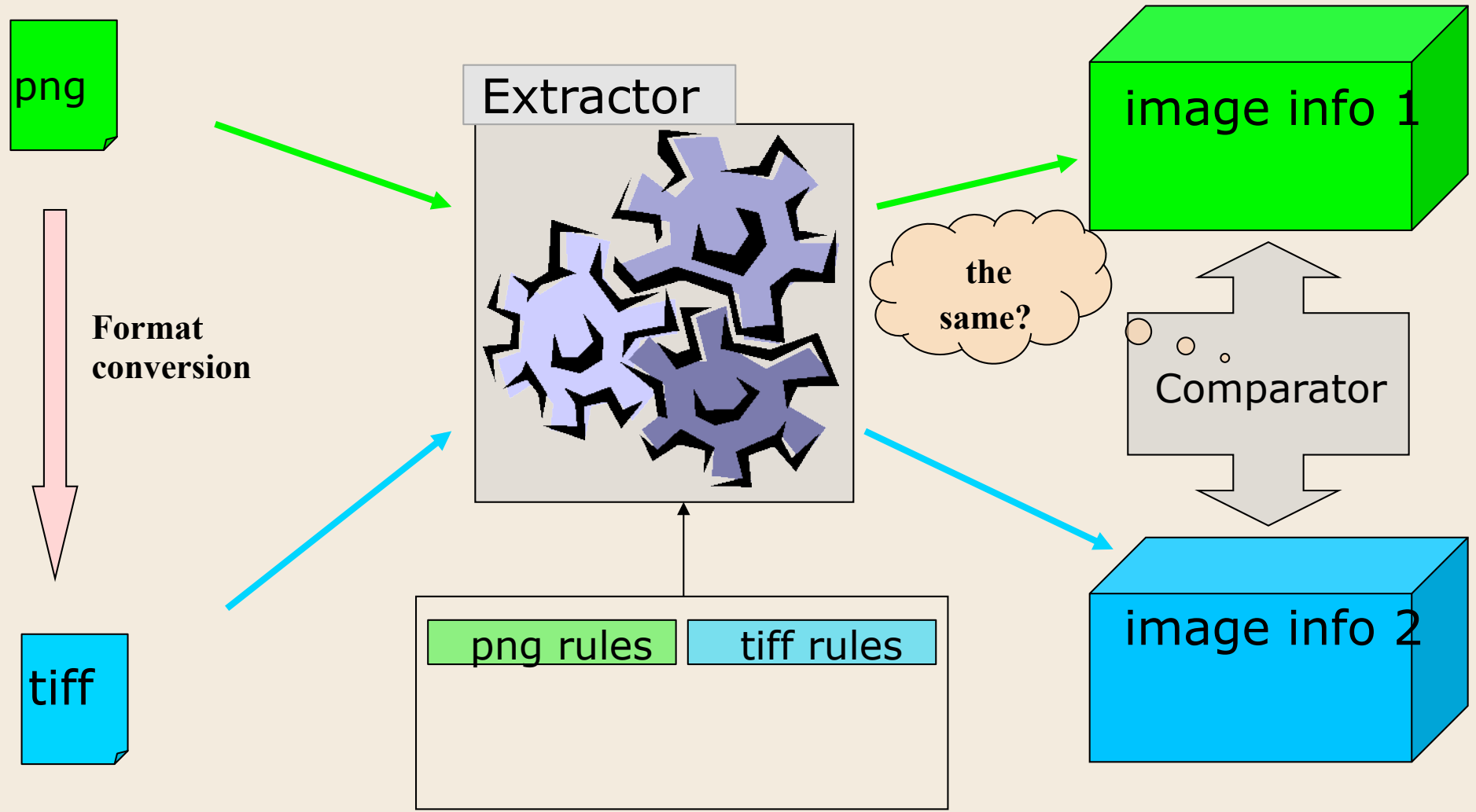


► Photoshop ►

# Planets: the vision 2

Historisch
Kulturwissenschaftliche
Informationsverarbeitung

Obj 1

Format conversion

Obj 2

Extractor

rule set 1   rule set 2

the same?

object info 1

Comparator

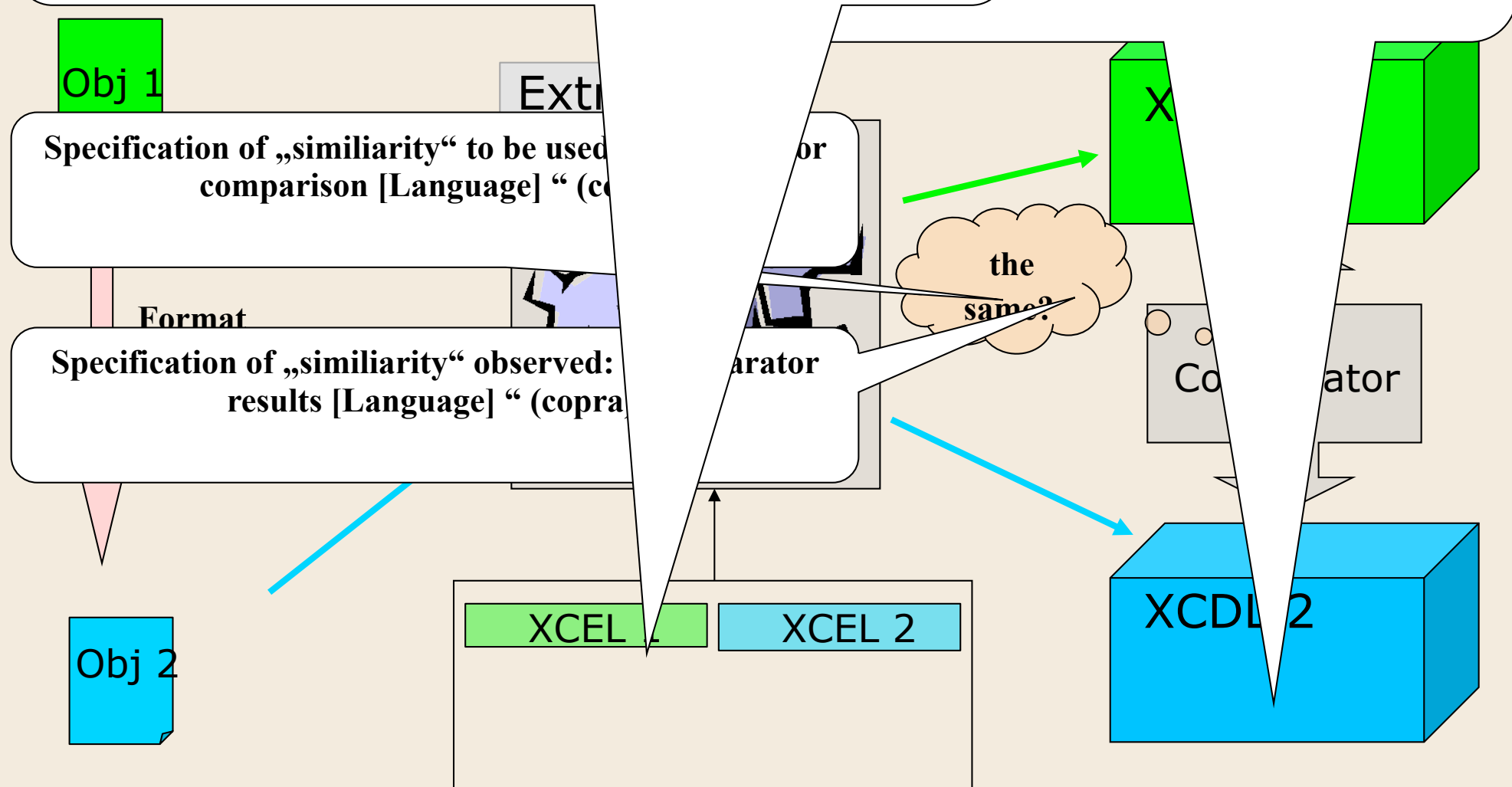object info 2

**Machine readable form of a file format specification: „eXtensible Characterisation Extraction Language" (XCEL), able to describe any machine readable format in a formal language, processible by a software tool for extraction of content as XCDL.**

...tion of file content: „eXtensible ...finition Language" (XCDL), able ...nt of digital objects (=1 + n more ...by a software tool for further analysis.

Obj 1

Extr...

X...

**Specification of „similiarity" to be used ...or comparison [Language] " (co...**

the same?

Format

Co...ator

**Specification of „similiarity" observed: ...arator results [Language] " (copra...**

Obj 2

XCEL 1   XCEL 2

XCDL 2

# Text in XCDL

This <sub>is a</sub> text

<refData id="1">**54 68 69 73** 20 69 73 20 61 20 <span style="color:blue">**74 65 78 74**</span></refData>
…
<property>
<name>fontsize</name>
<rawVal>
<val>48</val>
<type>unsignedInt8</type>
</rawVal>
<dataRef> <!-- property refers to discrete part of reference data-->
**<ref id="1" start="0" end="3"/>**
<span style="color:blue">**<ref id="1" start="10" end="12"/>**</span>
</dataRef>
</property>

# Image in XCDL

<refData id="1">7A 11 9B 77 34 89 72 11 29 F4 DA 9C B2 23 56 93 86 83 82 65 …</refData>
…
<property>
<name>title</name>
<rawVal>
<val>Ebstorf Mappa Mundi</val>
<type>ASCII</type>
</rawVal>
<dataRef>
<ref id="1" start="0" end="13455"/>
</dataRef>
</property>
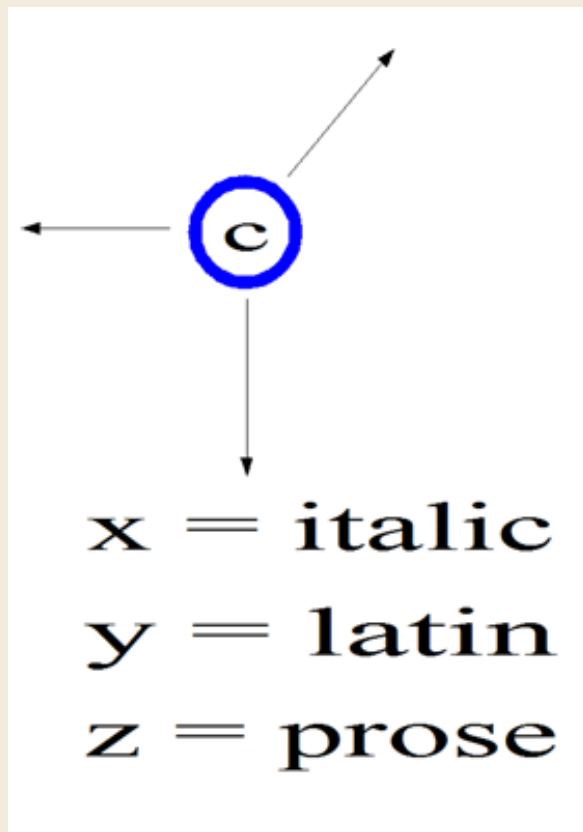
# Generalizing the practical solution

# Dimensions: geometry



x = 3 cm

y = 3 cm

z = 3 cm

Allows to make statements about the proximity of two objects on the "y" axis.

*Irrespective of the "shape" of the object.*

# Dimensions: textual / conceptual



x = italic
y = latin
z = prose

Allows to make statements about the proximity of two objects on the "y" axis.

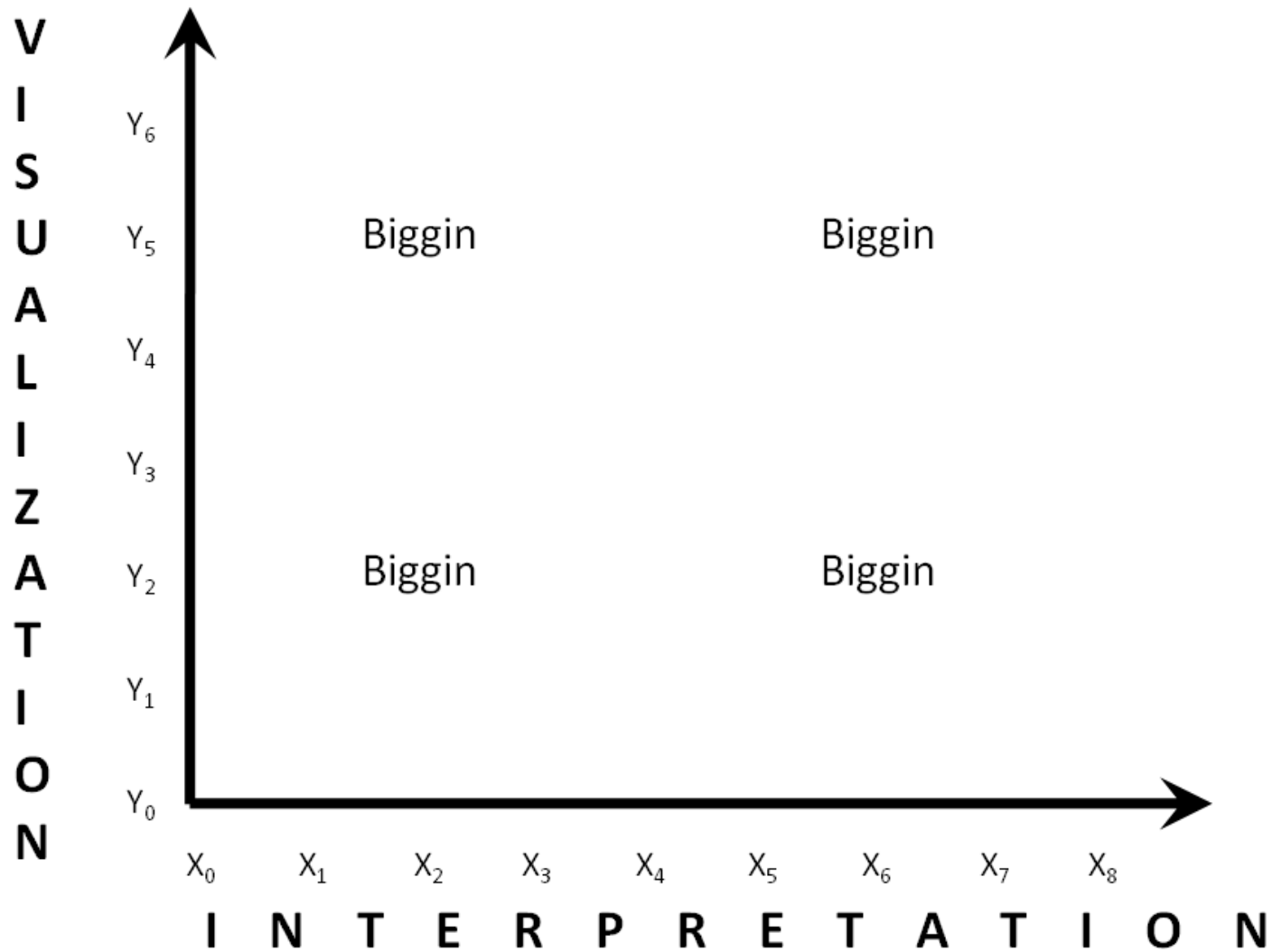Irrespective of the "object" that is at the abstract position.

# Dimensions: metrics
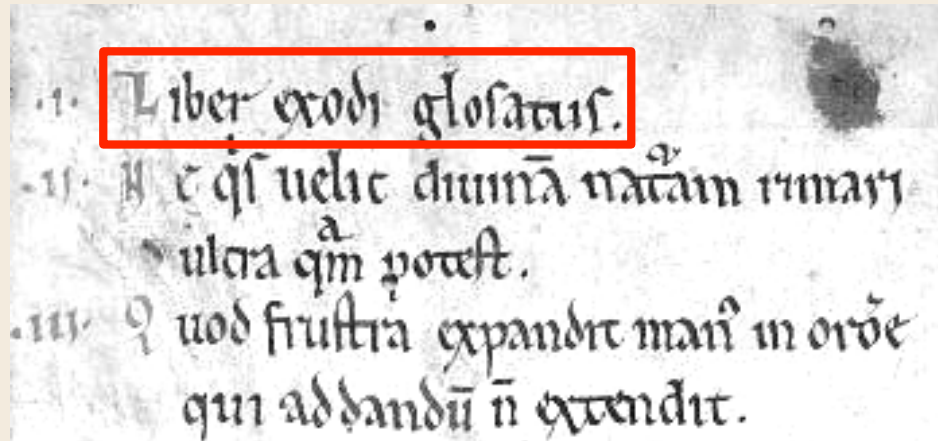
Dimensions are by definition orthogonal.

Dimensions can have any sort of metric:
- ❖ Rational: { - ∞ … + ∞ }
- ❖ Integer range: { 0 … 100 }
- ❖ Nominal: { medieval, early modern, modern }
- ❖ Image: {  ,  }
- ❖ …

# Four texts …

(1) &lt;person&gt;&lt;surname&gt;&lt;bold&gt;Biggin&lt;/bold&gt;&lt;/surname&gt;&lt;/person&gt;

(2) &lt;person&gt;&lt;surname&gt;&lt;italics&gt;Biggin&lt;/italics&gt;&lt;/surname&gt;&lt;/person&gt;

(3) &lt;airfield&gt;&lt;name&gt;&lt;bold&gt;Biggin&lt;/bold&gt;&lt;/name&gt;&lt;/airfield&gt;

(4) &lt;airfield&gt;&lt;name&gt;&lt;italics&gt;Biggin&lt;/italics&gt;&lt;/name&gt;&lt;/airfield&gt;

Which of the chunks are more similar to each other: (1) and (2) or (1) and (3)?
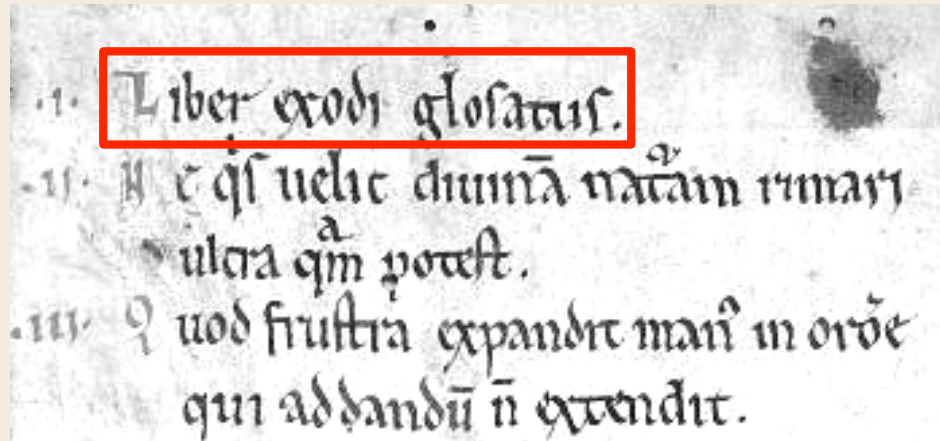
# … in a coordinate space.

# An image in a textual coordinate space



Liber exodi glosatus

# An text in an image coordinate space



Liber exodi glosatus

# An image in a semantic coordinate space



B.1 Bildbetrachter mit Darstellung der Polygone
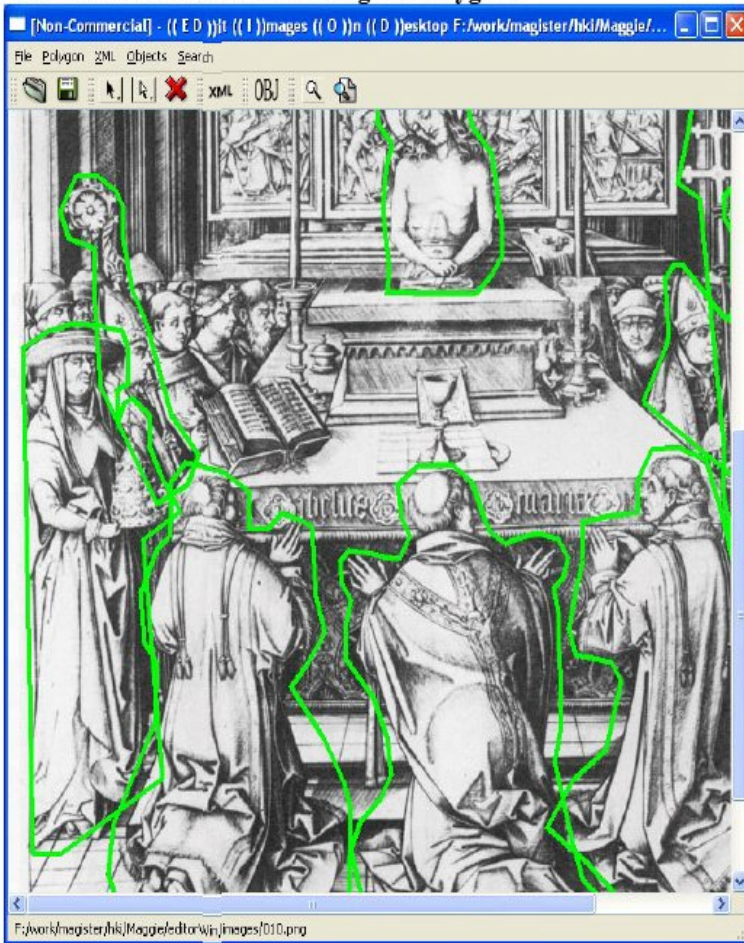
**Bishop**

**Cardinal**

**Monk**

**Priest**

**Monk**

# Semantics in an image coordinate space



B.1 Bildbetrachter mit Darstellung der Polygone
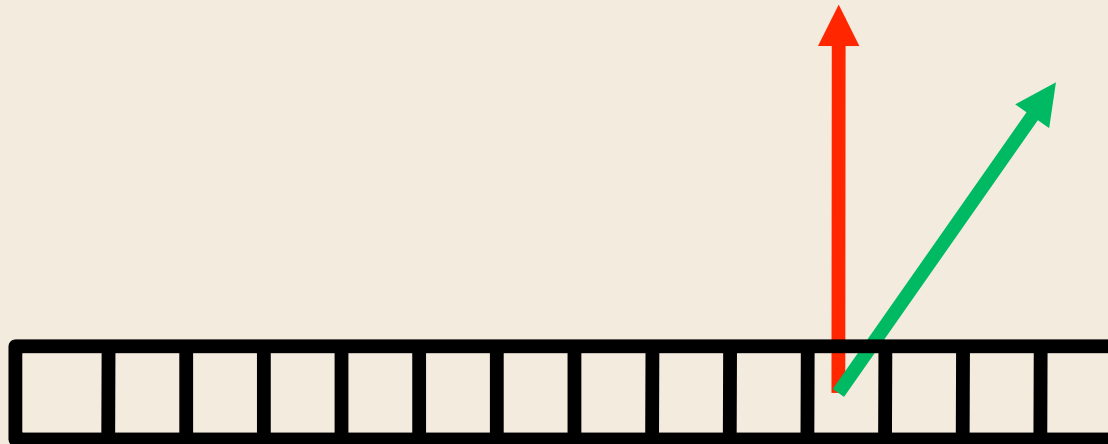
**Bishop**

**Cardinal**

**Monk**

**Priest**

**Monk**

# Generalization 1

**Biggin**

**<span style="color:red">Visualization {bold, italic}</span>**

**<span style="color:green">Interpretation {surname, topographic name}</span>**
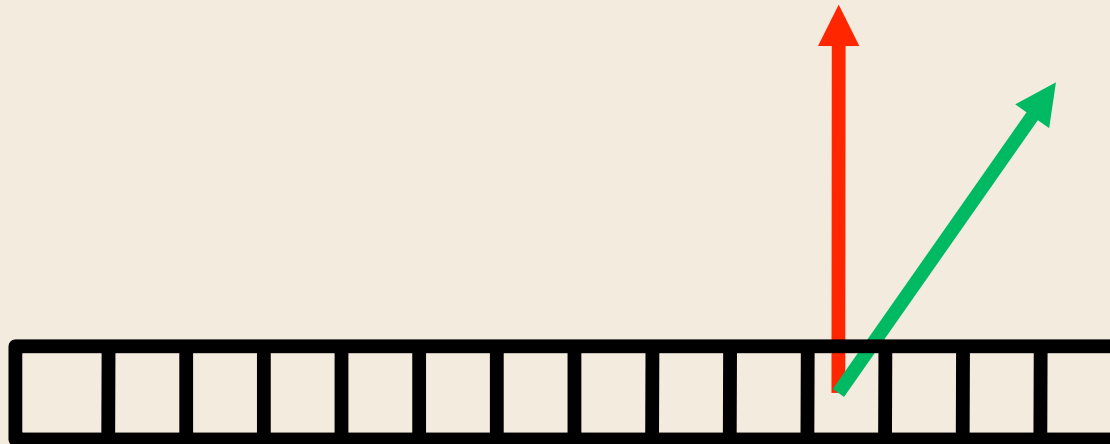
# Generalization 2



**Series of atomic content tokens**
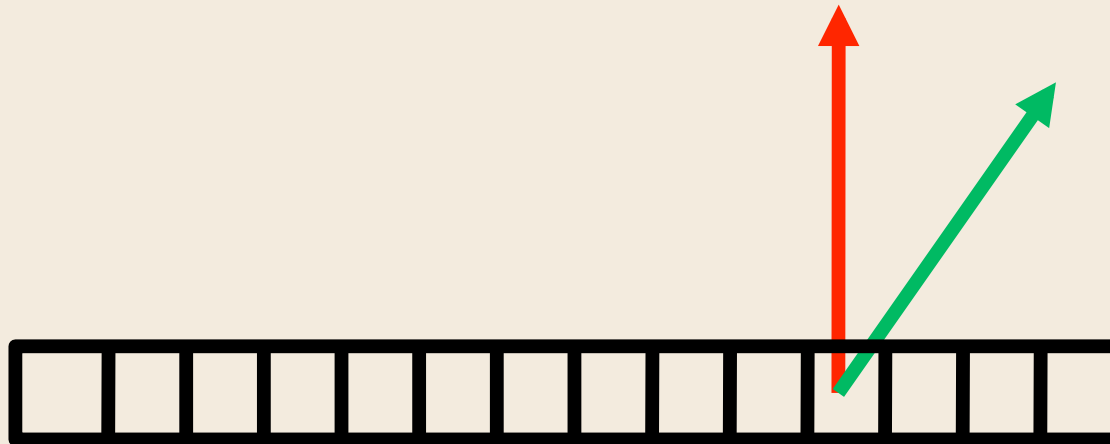**Conceptual dimension 1**
**Conceptual dimension 2**
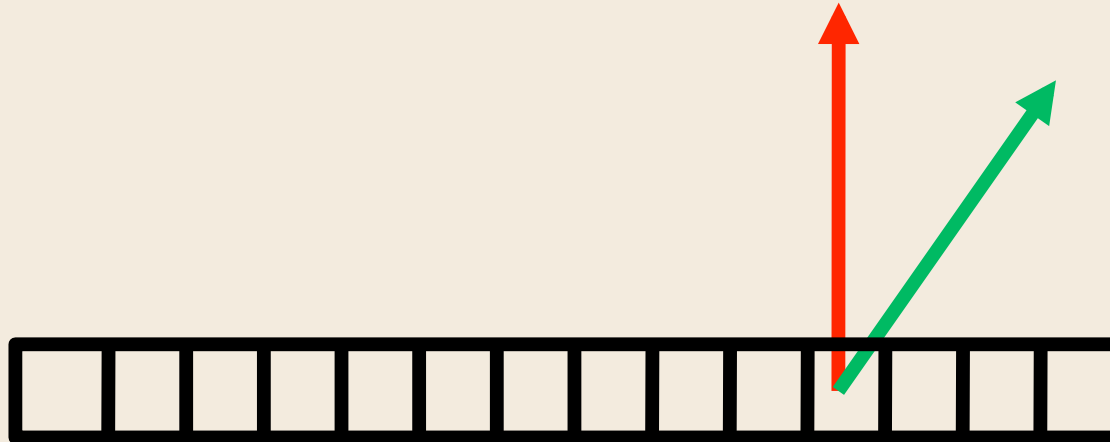
# Generalization 3

$\{ T, C_1, C_2\}$

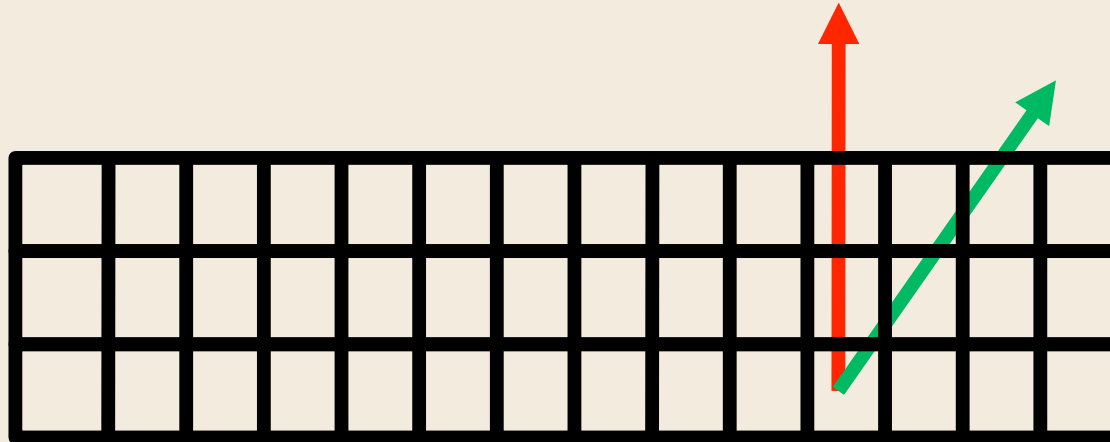# Generalization 4

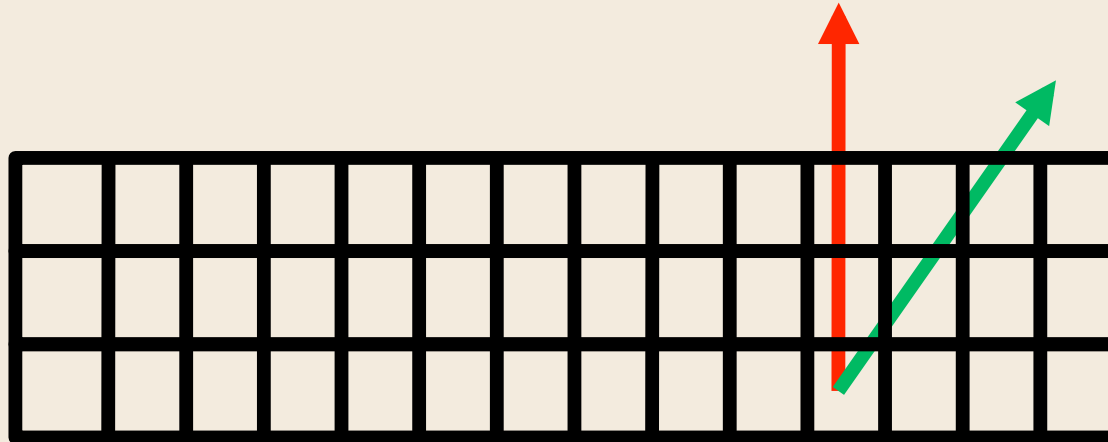$$\{ T, \{ C_1, C_2, \ldots, C_n \} \}$$

# Generalization 5

$\{ T, C_n \}$

(1) Texts are sequences of content carrying atomic tokens.

(2) Each of these tokens has a position in an *n*-dimensional conceptual universe.

# Generalization 6

$$\{ X, Y, C_n \}$$

# Generalization 7

$\{\ T_1,\ T_2,\ C_n\ \}$

(1) Images are planes of content carrying atomic tokens.

(2) Each of these tokens has a position in an $n$-dimensional conceptual universe.

# Generalization 8

$$I ::= \{ \{ T_1, T_2, \ldots T_m\}, C_n \}$$

(1) Information objects are $m$-dimensional arrangements of content carrying atomic tokens.

(2) Each of these tokens has a position in an $n$-dimensional conceptual universe.

# Generalization 9

$$I ::= \{T_m, \ C_n \}$$

(1) Information objects are $m$-dimensional arrangements of content carrying atomic tokens.

(2) Each of these tokens has a position in an $n$-dimensional conceptual universe.

(3) All of this, of course, is recursive …

# Another practical interlude
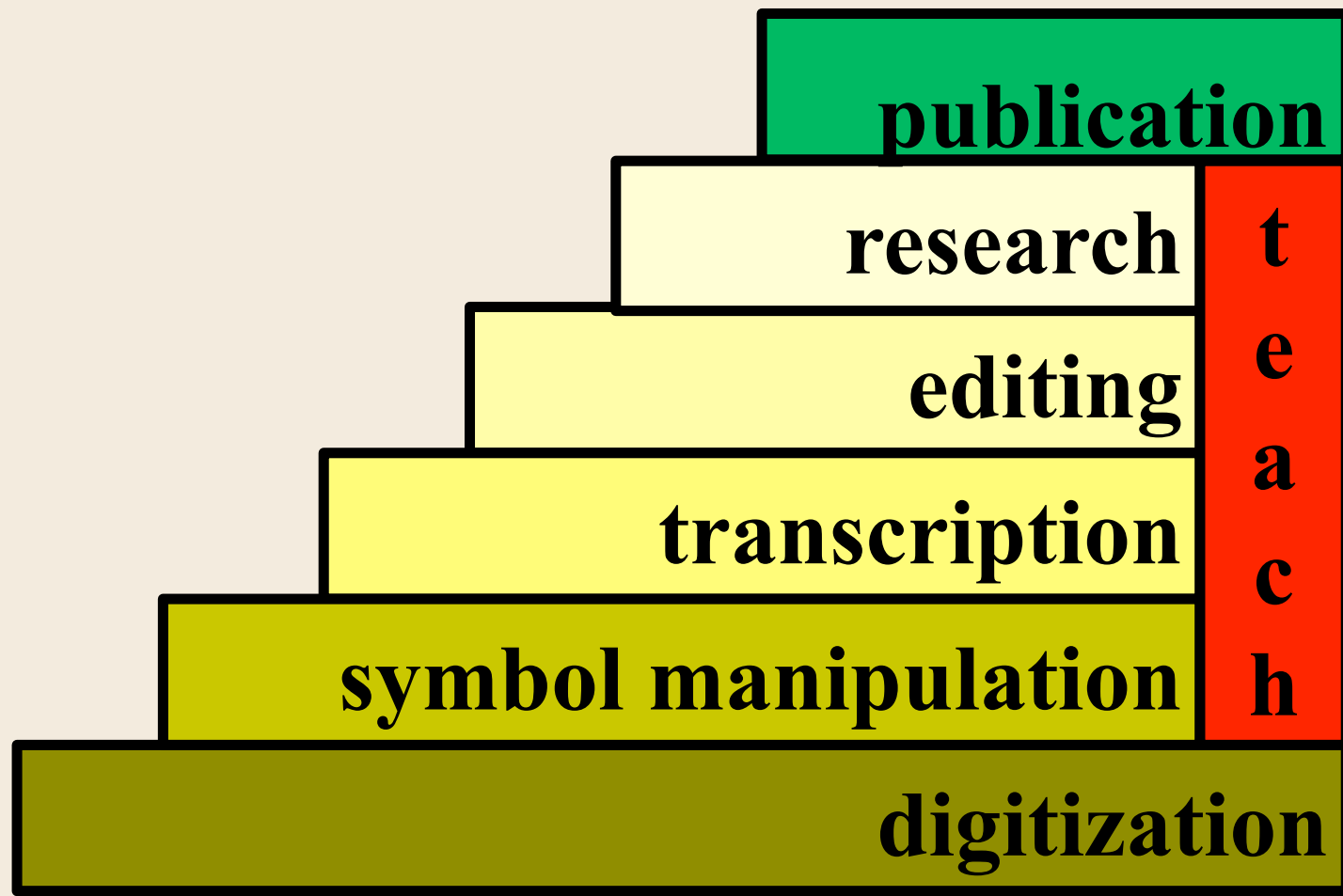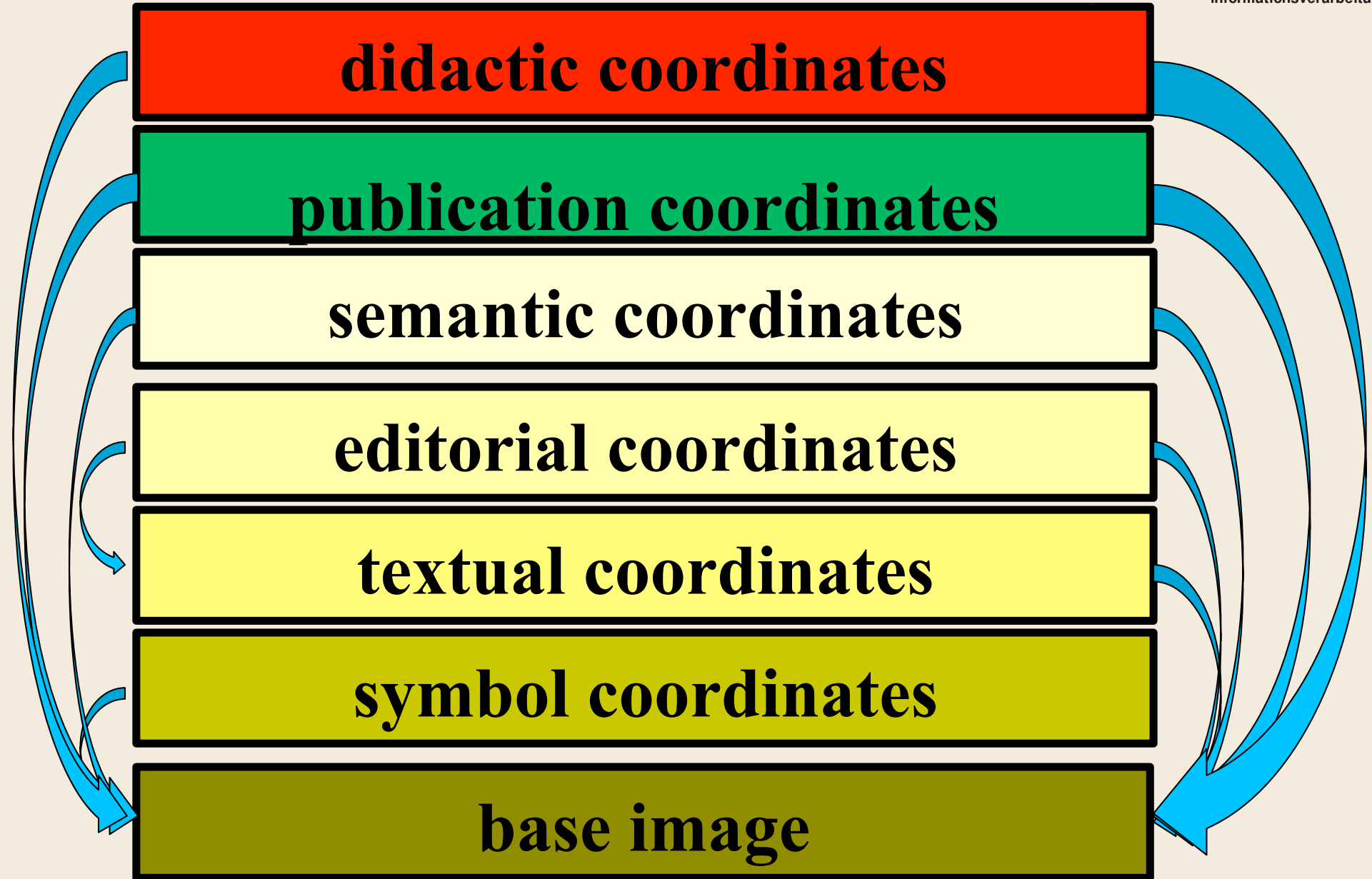
# Virtual Research Environments

http://www.monasterium.net/

→

Virtuelles deutsches Urkundennetz

(Virtual network of German charters)

# A model of historical research

publication

research

editing

transcription

symbol manipulation

digitization

teach

# A model of historical research

**didactic coordinates**

**publication coordinates**

**semantic coordinates**

**editorial coordinates**

**textual coordinates**

**symbol coordinates**

**base image**

# Conclusion

# Summary

(1) All texts, for which we cannot consult the producer, should be understood as a sequence of tokens, where we should keep the representation of the tokens and the representation of our interpretation thereof completely separate.

(2) Such representations can be grounded in information theory.

(3) These representations are useful as blueprints for software on highly divergent levels of abstraction.

# Thank you!

**manfred.thaller@uni-koeln.de**