

# **Evaluation in Computational Creativity**

Hannu Toivonen University of Helsinki www.cs.helsinki.fi/hannu.toivonen

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

28.9.2015 165



- Evaluation of creativity allows us to compare methods, to control progress, to improve methods
- However, evaluation of creativity is very difficult
  - No precise definition of creativity
  - Various goals (novelty, value, originality, ...)
  - Context-dependence
  - Cost of evaluation

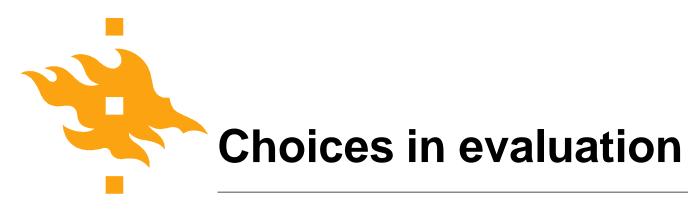
— ...



The goal of evaluation should be aligned with the goals of the system. E.g.:

- Machine creativity:
  Creative performance of creative programs
- Computer-supported creativity: Increase in creativity of humans using CC tools
- Creativity studies: Increase in knowledge about creative processes

#### - Focus here: evaluation of machine creativity



- 1. Summative vs. formative evaluation
- Is the goal to rate or compare systems (summative evaluation) or to help develop them (formative evaluation)?
- 2. Expert vs. layman vs. peer evaluation
- Who carries out the evaluation? The developers know the internals of their system best, laymen can provide objective views of the results. Peers have been an under-used evaluation resource in CC.
- 3. Evaluation of outcome vs. process (next slide)



## **Evaluation of Machine Creativity**

Two possible targets in evalution of machine creativity (Colton 2008):

- Artefact-based evaluation: are the results creative?
  - e.g: novelty and value of results
- Process-based evaluation: is the process creative?
  - e.g: combinatorial/ exploratory/ transformational creativity; creative acts of the FACE model



## Ritchie's Framework for Artefact Based Evaluation

Ritchie (2007)

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI



Consider a set R of artefacts produced by a system. Primitive properties that can be considered:

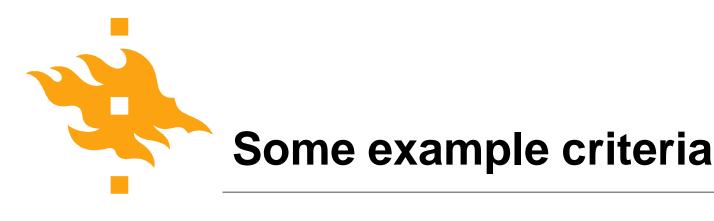
- Typicality: Is the artefact a typical/ recognizable example of the target genre?
- Novelty: How (dis)similar is the artefact to existing examples of its genre?
- Quality [= Value]



- typ(a) = amount of typicality associated to artefact a
- val(a) = amount of quality associated to a
- $\mathsf{T}_{\alpha,\beta}(\mathsf{X}) = \{ a \in \mathsf{X} \mid \alpha \leq \mathsf{typ}(a) \leq \beta \}$ 
  - Set of artefacts a with typicality between  $\alpha$  and  $\beta$

$$- V_{\alpha,\beta}(X) = \{ a \in X \mid \alpha \leq val(a) \leq \beta \}$$

- Set of artefacts a with value between  $\alpha$  and  $\beta$
- size(X) = number of elements of X
- ratio(X,Y) = size(X) / size(Y)
- R: a set of artefacts produced by the system



 $\begin{array}{ll} \underline{Criterion\ 2} & ratio(\mathsf{T}_{\alpha,1}(\mathsf{R}),\,\mathsf{R}) > \theta \\ & - \mbox{ at least fraction } \theta \mbox{ of results } \mathsf{R} \mbox{ have high typicality } (>\alpha) \\ \underline{Criterion\ 4} & ratio(\mathsf{V}_{\gamma,1}(\mathsf{R}),\,\mathsf{R}) > \theta \\ & - \mbox{ at least fraction } \theta \mbox{ of results } \mathsf{R} \mbox{ have high value } (>\gamma) \\ \underline{Criterion\ 5} & ratio(\mathsf{V}_{\gamma,1}(\mathsf{R})\ \cap\ \mathsf{T}_{\alpha,1}(\mathsf{R}),\,\mathsf{T}_{\alpha,1}(\mathsf{R})) > \theta \\ & - \mbox{ at least fraction } \theta \mbox{ of typical } (>\alpha) \mbox{ results also have high value } (>\gamma) \end{array}$ 



- Any creative system is based on some existing examples, in one way or another. These can – and should – be taken into account.
- The *inspiring set I* consists of all the relevant artefacts known to the program designer, or items which the program is designed to replicate, or a knowledge base of known examples which drives the computation within the program
- Inspiring set ≈ training set in ML/DM

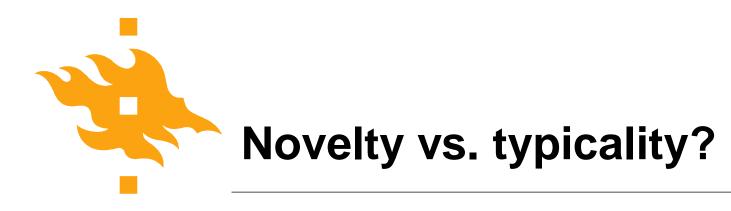


### <u>Criterion 9</u> ratio( $I \cap R, I$ ) > $\theta$

- Results R reproduce at least fraction  $\theta$  of the inspiring set I
- Is the system able to reproduce its inspiring set?
- Cf. ML: are training example classified correctly?

### <u>Criterion 10</u> ratio(R, $I \cap R$ ) > $\theta$

- Results R contain at least θ-1 times as many items outside the inspiring set I as inside it
- Can the system extrapolate/generalize outside the inspiring set/training examples?



Novelty and typicality are subtly different:

- Not recognizable as a member of the genre
   → low typicality
- Very different from the inspiring set (but possibly very clearly within the genre)
  - $\rightarrow$  high novelty



Note: Ritchie does not prescribe a set of criteria. Instead, the criteria must be designed and chosen according to the goals and needs of each work; Richie gives examples of some of the possible criteria that one may want to use .

## FACE Model for Process-Based Evaluation

Pease and Colton (2011)

HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

28.9.2015 178



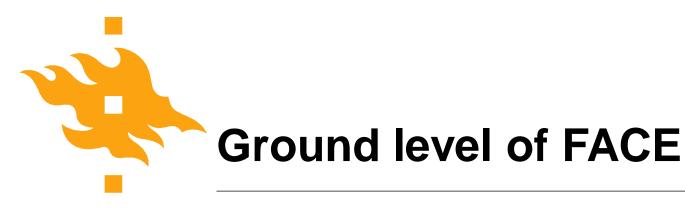
- Focus on creative processes, not their results
- In the FACE model, systems can be characterized by their creative acts
- The four aspects of the model:
  - F framing
  - A aesthetics
  - C concept
  - E expression
- Here we present a simplified version



- C: the concept or the idea of the artefact
  - E.g. use of excessive rhyming in poetry
- E: a concrete expression of the concept
  - E.g. a poem that uses excessive rhyming
- A: a measure of *aesthetics* of the work of art
  - E.g. grammaticality etc. of a poem
- F: all background information about the piece (*framing*)
  - E.g. a description of why excessive rhyming could be interesting, and what the poem expresses



- Framing is especially important for computational creativity
- It is difficult to appreciate the output (expression) without knowing anything about the process, its goals, etc.
  - E.g., is the resulting image pretty just by change? Or did the system produce it based on some specific criteria and goals? Was the process complicated? Is there some intention, e.g., a message that is being conveyed?



- Ground-level generative acts and their products
  - Act  $F^g \rightarrow$  generates an item of framing information
  - Act  $A^g \rightarrow$  generates an aesthetic measure
  - Act  $C^g \rightarrow$  generates a concept
  - Act  $E^g \rightarrow$  generates an expression of a concept
- Any system can now be described in terms of who carries out these acts, and how
  - A simple generative system only performs E<sup>g</sup>
  - A system that learns to evaluate also performs A<sup>g</sup>
  - (The programmer and other humans probably perform the other acts)



- FACE also has a meta-level: processes that produce ground-level generators
- Process-level acts and their outputs:
  - Act  $F^p \rightarrow$  generates a method for generating framing information
  - Act  $A^p \rightarrow$  generates a method for generating aesthetic measures
  - Act  $C^p \rightarrow$  generates a method for generating concepts
  - Act E<sup>p</sup> → generates a method for generating expressions of a concept







HELSINGIN YLIOPISTO HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI

www.helsinki.fi/yliopisto



- $F^p$ : Methods for generating the contextual history of this genre of art
- $F^g$ : The contextual history of this genre of art, motivation, justification, etc.
- $A^p$ : Methods for generating the idea of art having multiple meanings when viewing from multiple perspectives
- A<sup>g</sup>: The idea of art having multiple meanings when viewing from multiple perspectives
- $C^p$ : Methods for generating new perspectives from which the art might make sense
- $C^g$ : The constraint that a picture must make sense when upside down
- $E^p$ : Methods for generating expressions of art which have a different meaning when viewed upsidedown
- $E^g$ : Expressions of art which have a different meaning when viewed upsidedown (see figure 1)

HELSINGIN YLIOPI HELSINGFORS UNIVERSITET UNIVERSITY OF HELSINKI