
Minimum Description Length

Principle: between Theory and Practice

”Theory without practice is empty, practice without theory is blind”

Prof. Alexey Potapov

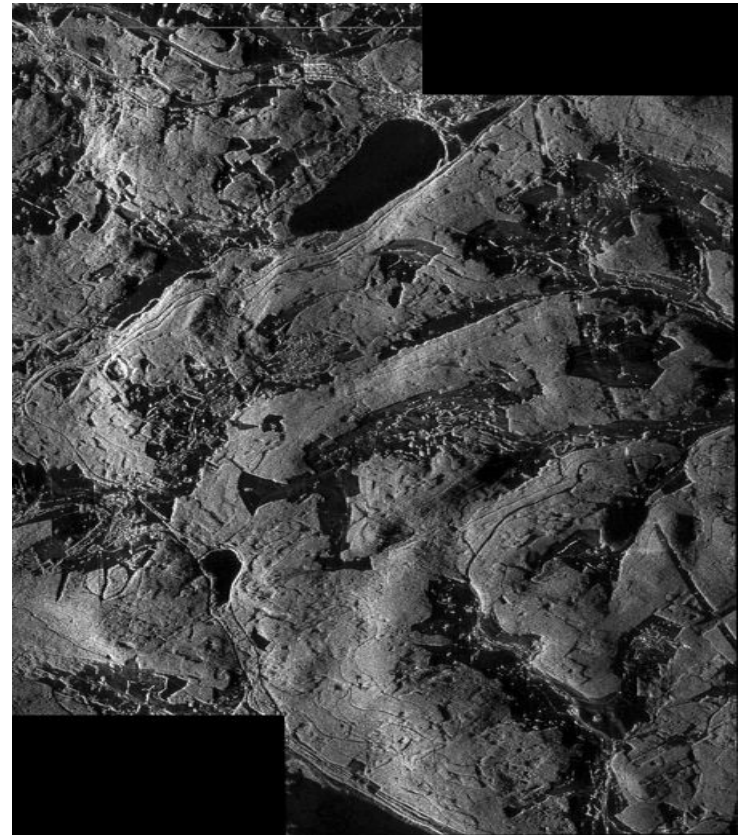
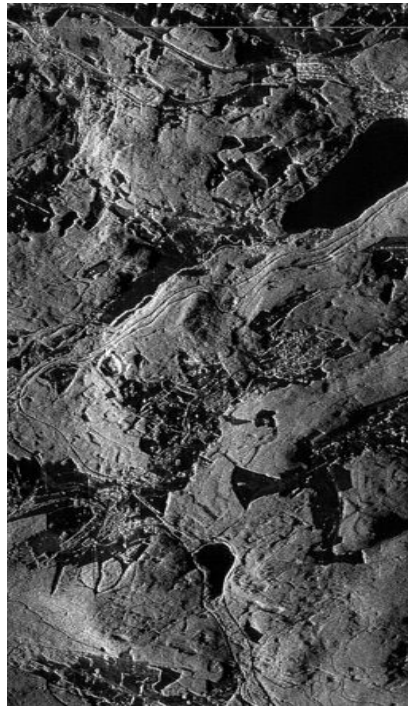
ITMO University, Saint-Petersburg State University, AIDEUS

2015

AGI'15 @ Berlin

One practical task: image matching

- How to find correspondence between pixels of two images of the same scene?



Simplest approach: correlation

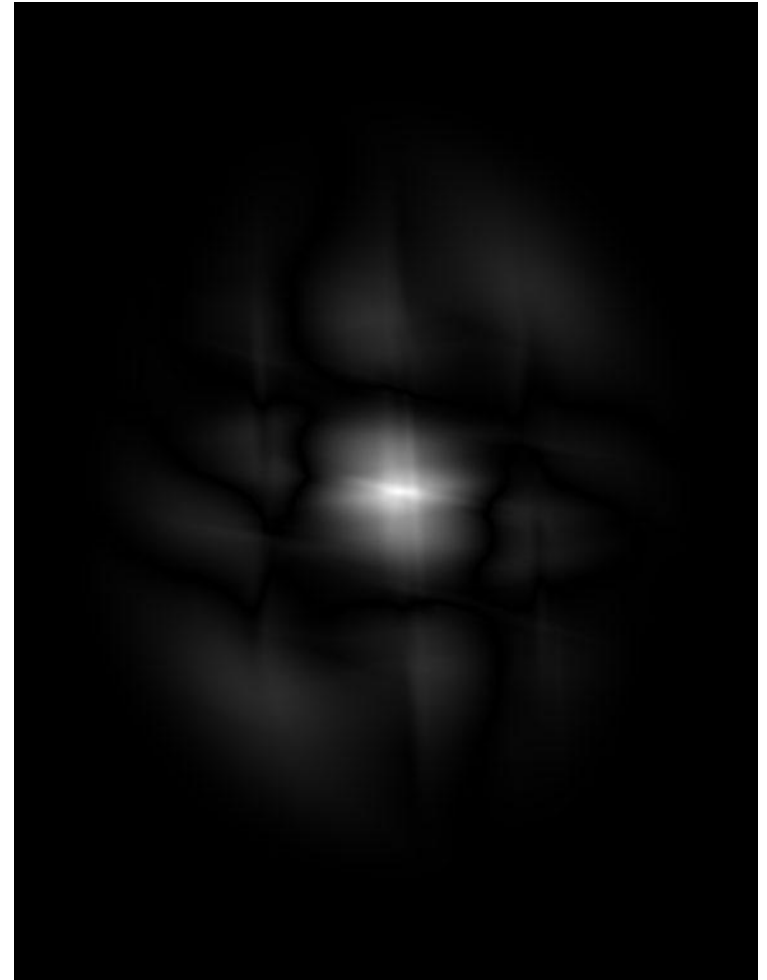
Least squares error

$$E_{f_1, f_2}(x, y) = \frac{1}{N^2} \sum_{x=0}^{N-1} \sum_{y=0}^{N-1} (f_1(x, y) - f_2(x, y))^2$$

Correlation

$$C_{f_1, f_2}(x, y) = \frac{1}{N^2} \sum_{x=0}^{N-1} \sum_{y=0}^{N-1} f_1(x, y) f_2(x, y)$$

Slightly more advanced:
cross-correlation function calculated
via Fourier Transform



Fourier-Mellin Transform

1. Amplitude spectrum
2. Log-polar transform
3. Cross-corr. Via Fourier
4. Find scale/rotation
5. Compensate scale/rotation
6. Cross.corr. to find shifts
7. Success!

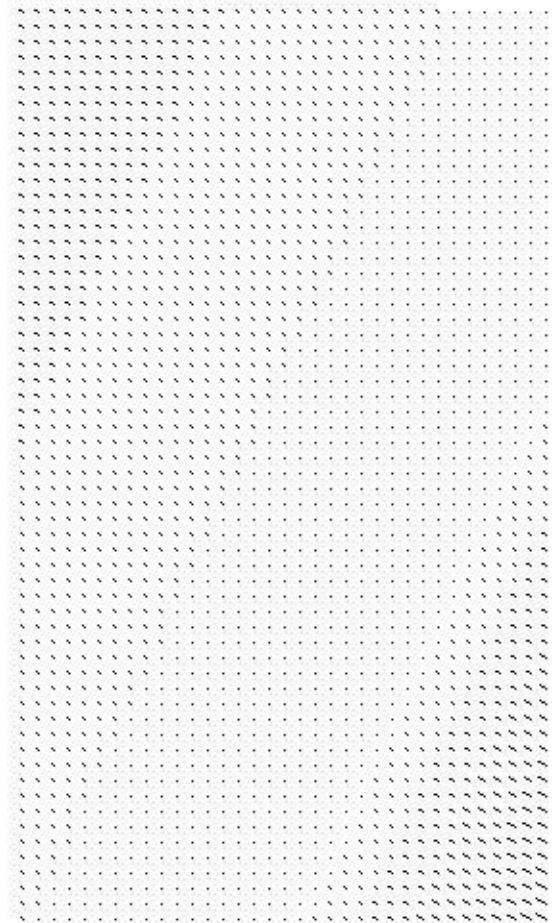


Block Matching: Local displacement extension

1. Take local fragments around different points of pre-aligned images
2. Match them by correlation
3. Construct local displacement field



Resulting displacement field



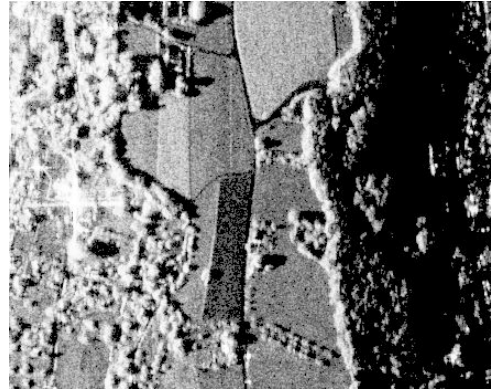
General solution for aerospace image matching!?

However...

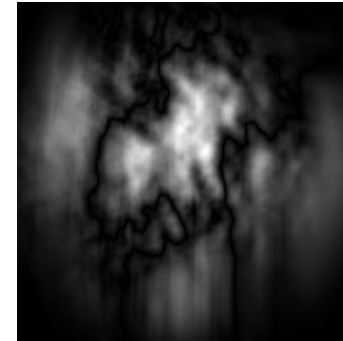
Optical image



SAR image



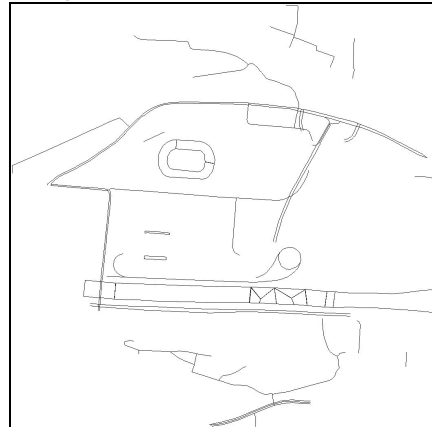
Cross-correlation field



Optical image



Digital map



Correlation?

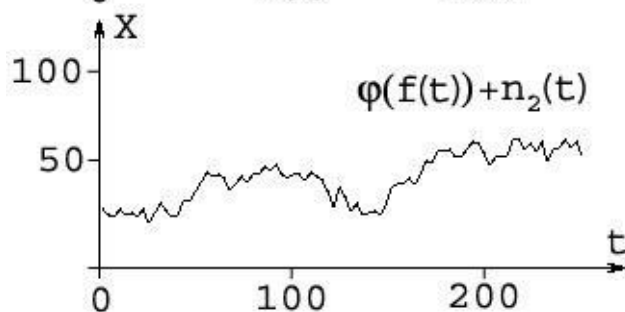
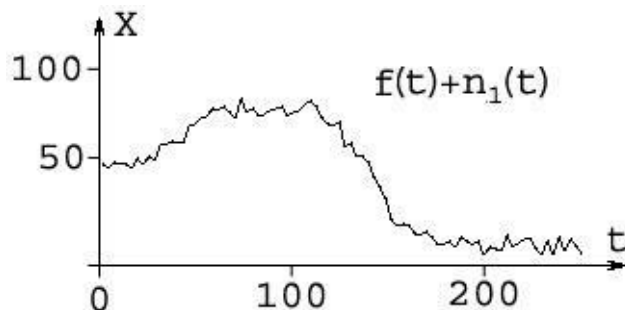
Many applications require matching images of different modalities

Criterion: Mutual Information

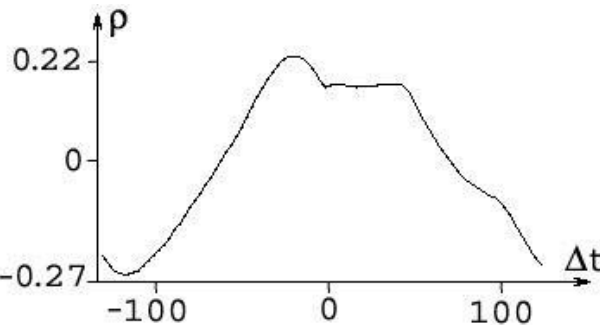
Mutual information $I(X, Y) = E_{XY} [I(x; y)] = \sum_{x \in X} \sum_{y \in Y} P(x, y) \log_2 \frac{P(x, y)}{P(x)P(y)}$

No correlation $\Rightarrow E_{XY} [XY] - E_X [X]E_Y [Y] = 0$

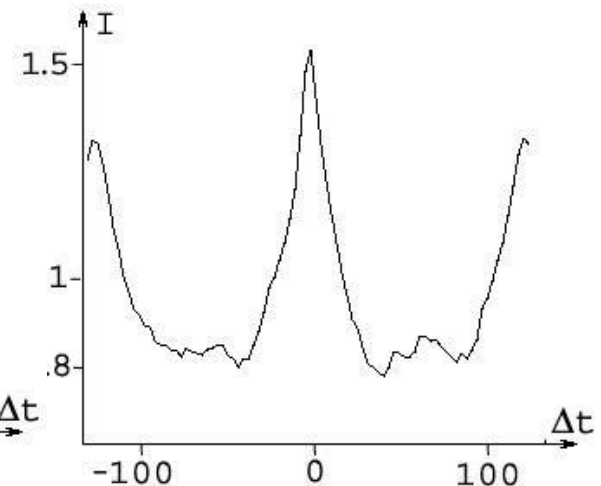
No mutual information $\Rightarrow E_{XY} [\varphi_1(X)\varphi_2(Y)] - E_X [\varphi_1(X)]E_Y [\varphi_2(Y)] = 0$



Cross correlation:
degraded maximum

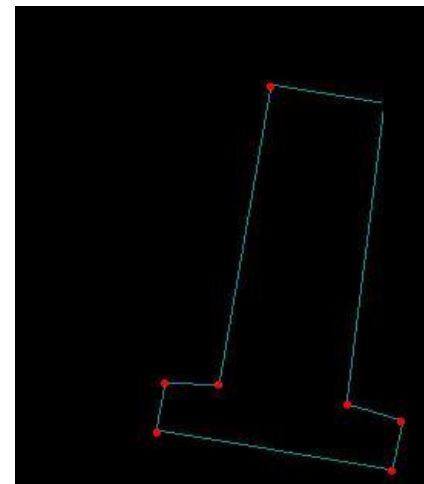
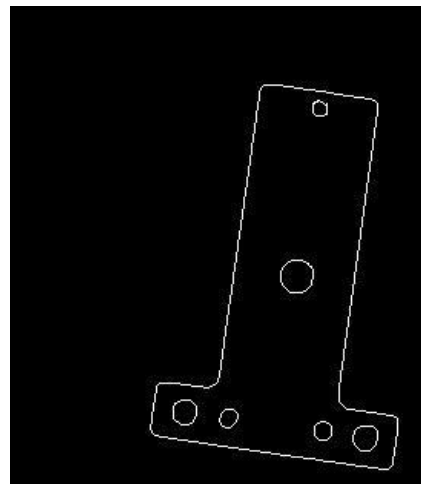
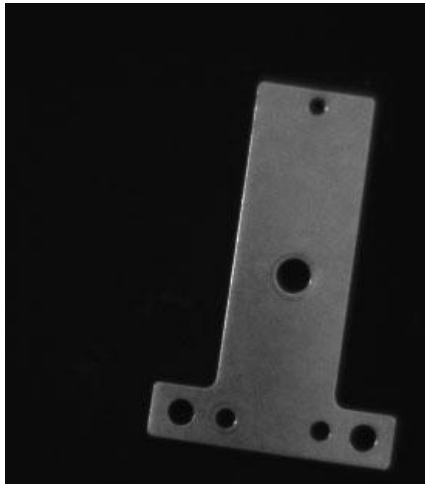
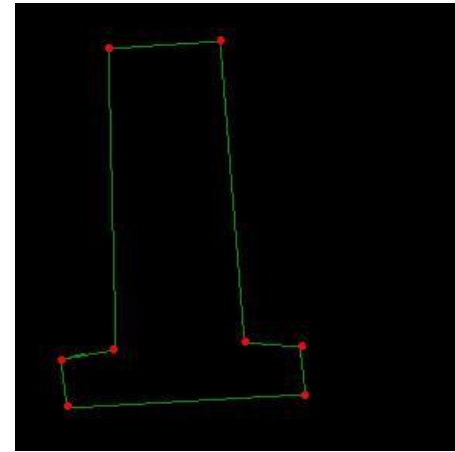
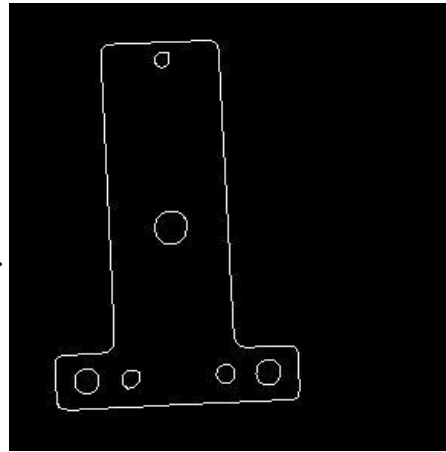
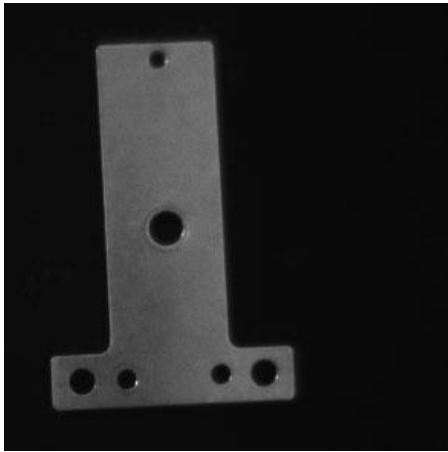


Mutual information:
Ideal maximum



Unfortunately, it's difficult to compute and not applicable to vector maps

Invariant structural descriptions



Image

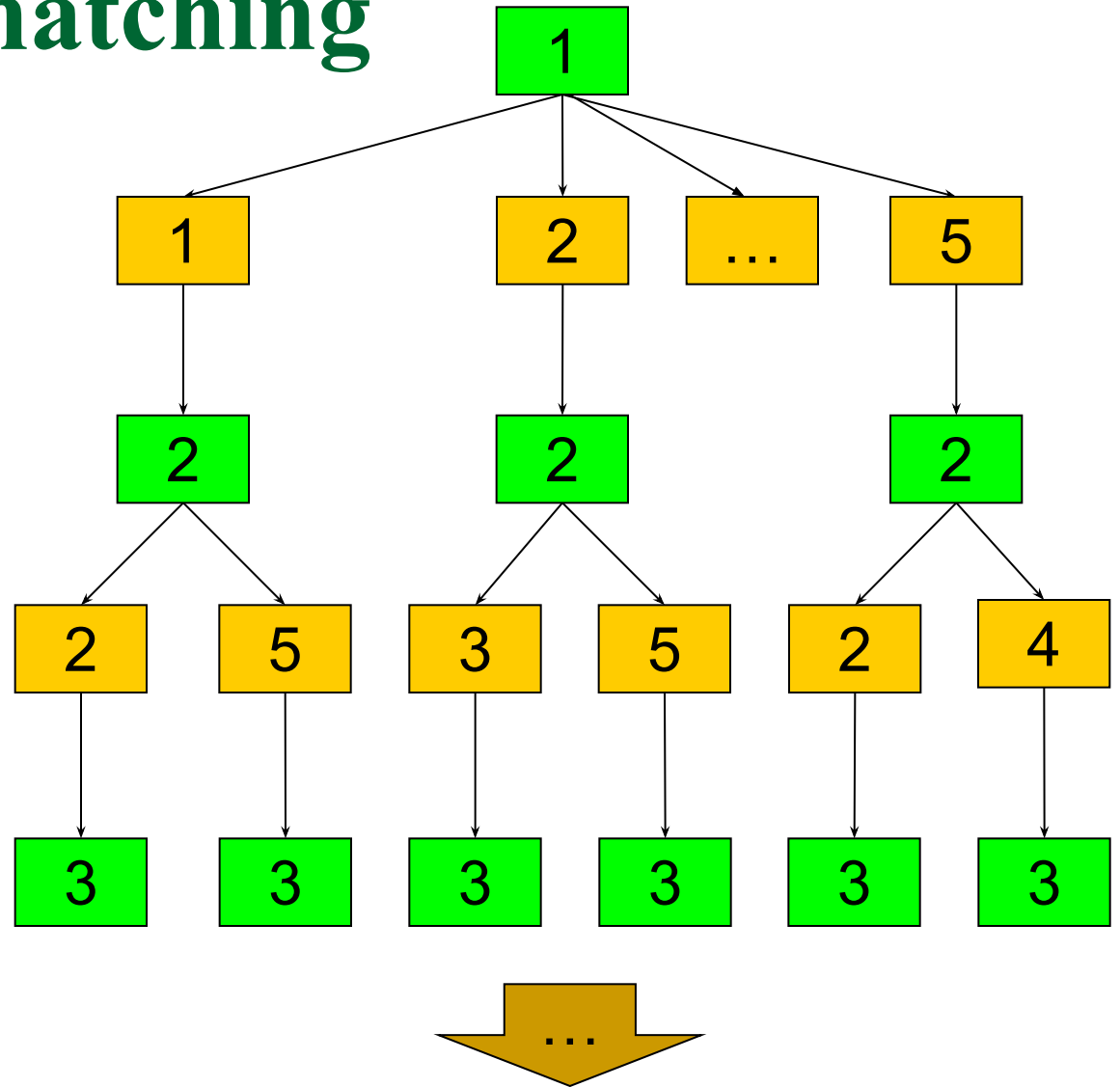
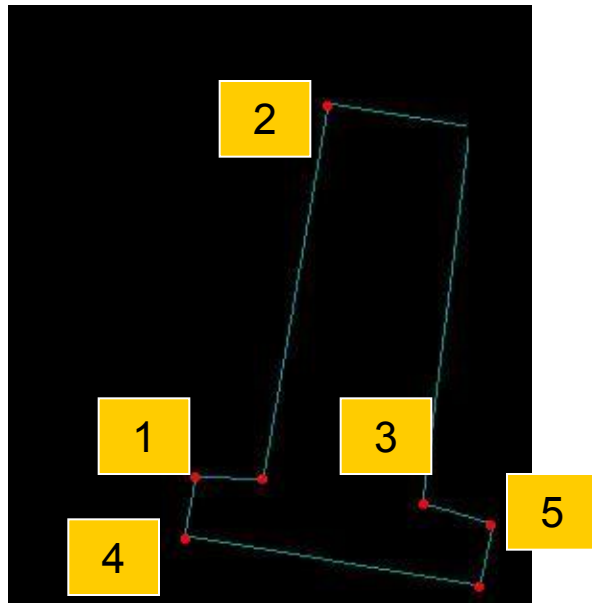
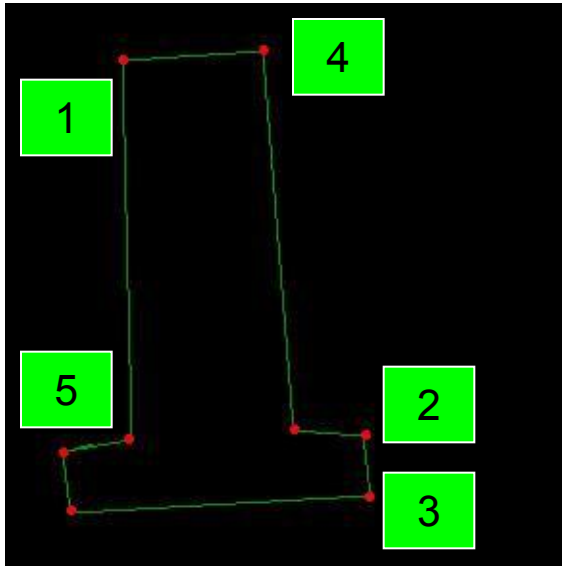


Contours



Structural elements

Structural matching

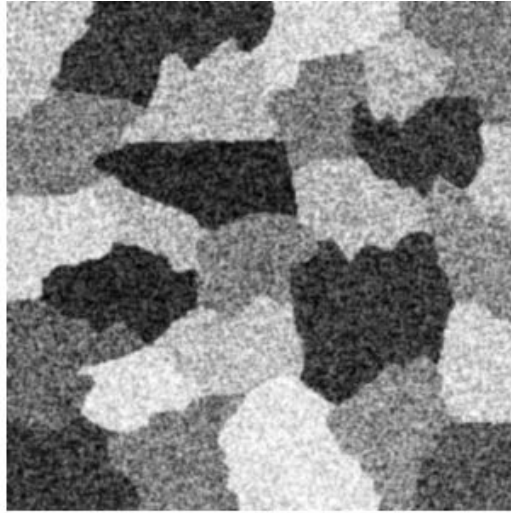


More questions...

- How to estimate quality of structural correspondence?
 - How to choose the group of transformations if it is not known?
- How to construct contours and structural elements optimally?
 - How to choose the most adequate number of contours and structural elements?
- Are precision criteria such as mean square error suitable? Or have they the same shortcomings as correlation?

MSE criterion: oversegmentation

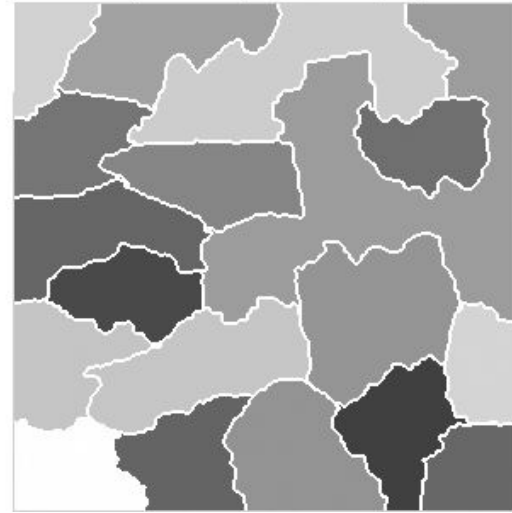
Each region is described by average value



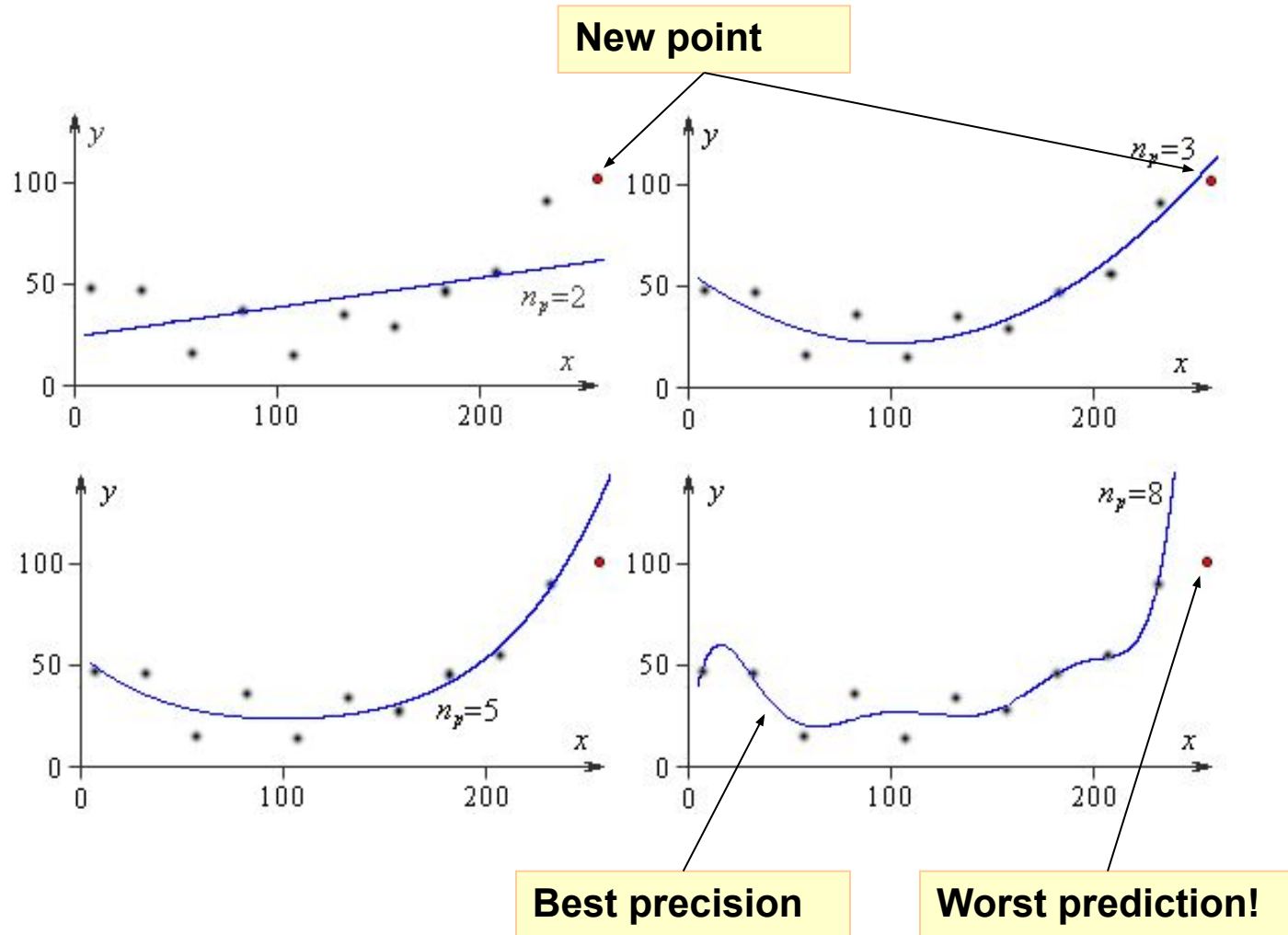
Correct, but not the most precise description!

More precise

Over-segmentation!



Functional approximation



Information-theoretic criterion

Again, criteria from information theory help:

- Mutual information can be extended for the task of matching structural elements
- In general, the minimum description length can be used for model selection

The best model is the model that minimizes the sum

- the description length (in bits) of the model,
- the description length (in bits) of data encrypted with help of the model (deviation of data from model).

Connection to Bayes' rule

$$\text{Bayes rule: } P(H | D) = \frac{P(D | H)P(H)}{P(D)}$$

- Posterior probability: $P(H | D)$
- Prior probability: $P(H)$
- Likelihood: $P(D | H)$

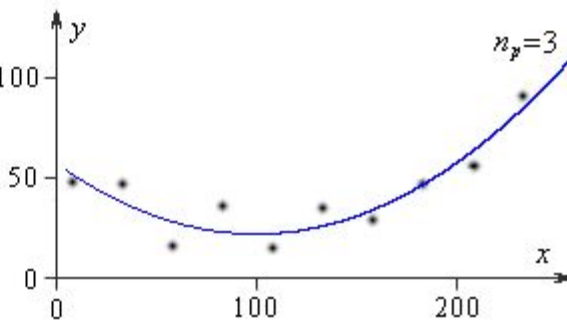
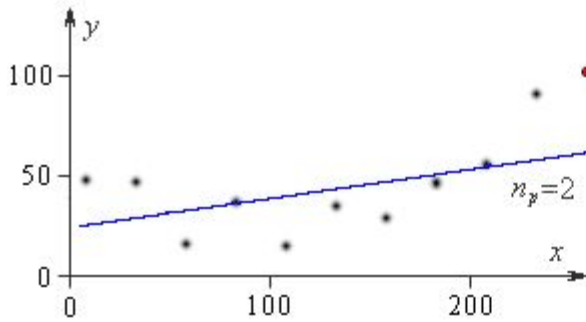
$$H^* \square \operatorname{argmax}_H P(H | D) \square \operatorname{argmax}_H P(H)P(D | H) \square \square$$

$$\square \operatorname{argmin}_H \square \log P(H) \square \log P(D | H) \square$$

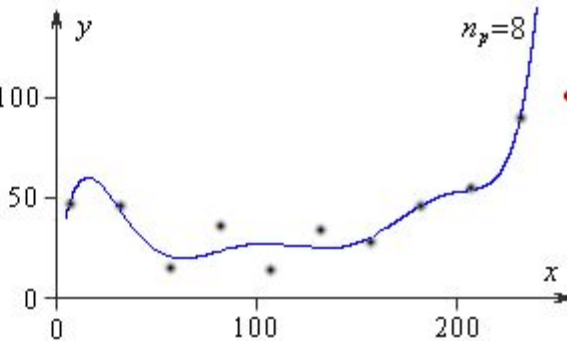
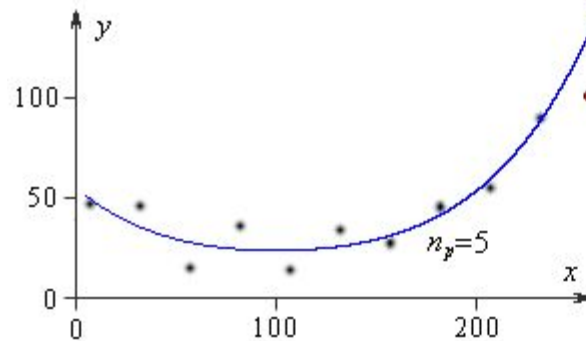
- The description length of the model: $-\log P(H)$
- The description length of data encrypted with the help of the model: $-\log P(D | H)$

Application to function approximation

Too simple model



The best model is chosen as trade-off between precision and complexity



Too complex model

$$L = \underbrace{\frac{n_p}{2} \log_2 n}_{l(H)} + \underbrace{\frac{n}{2} \log_2 \frac{\varepsilon^2(\mathbf{w})}{n}}_{K(D|H)}$$

$$K(D | H) = -\log_2 P(D | H) = -\log_2 \prod P(\varepsilon_i(\mathbf{w}))$$

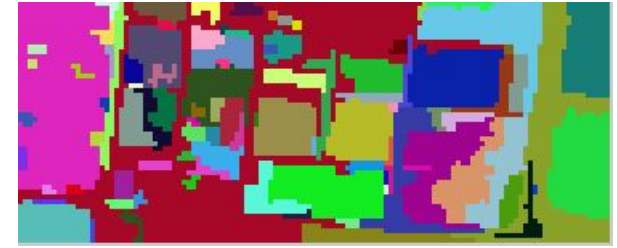
Application to image segmentation



Initial image



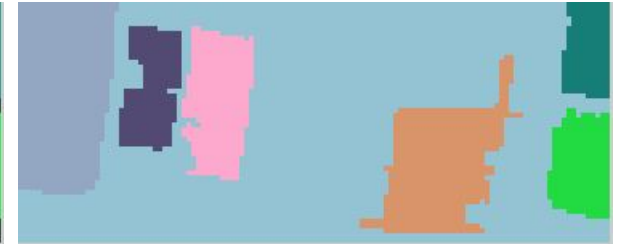
$N_{gr}=300$; $DL=4,5e+5$



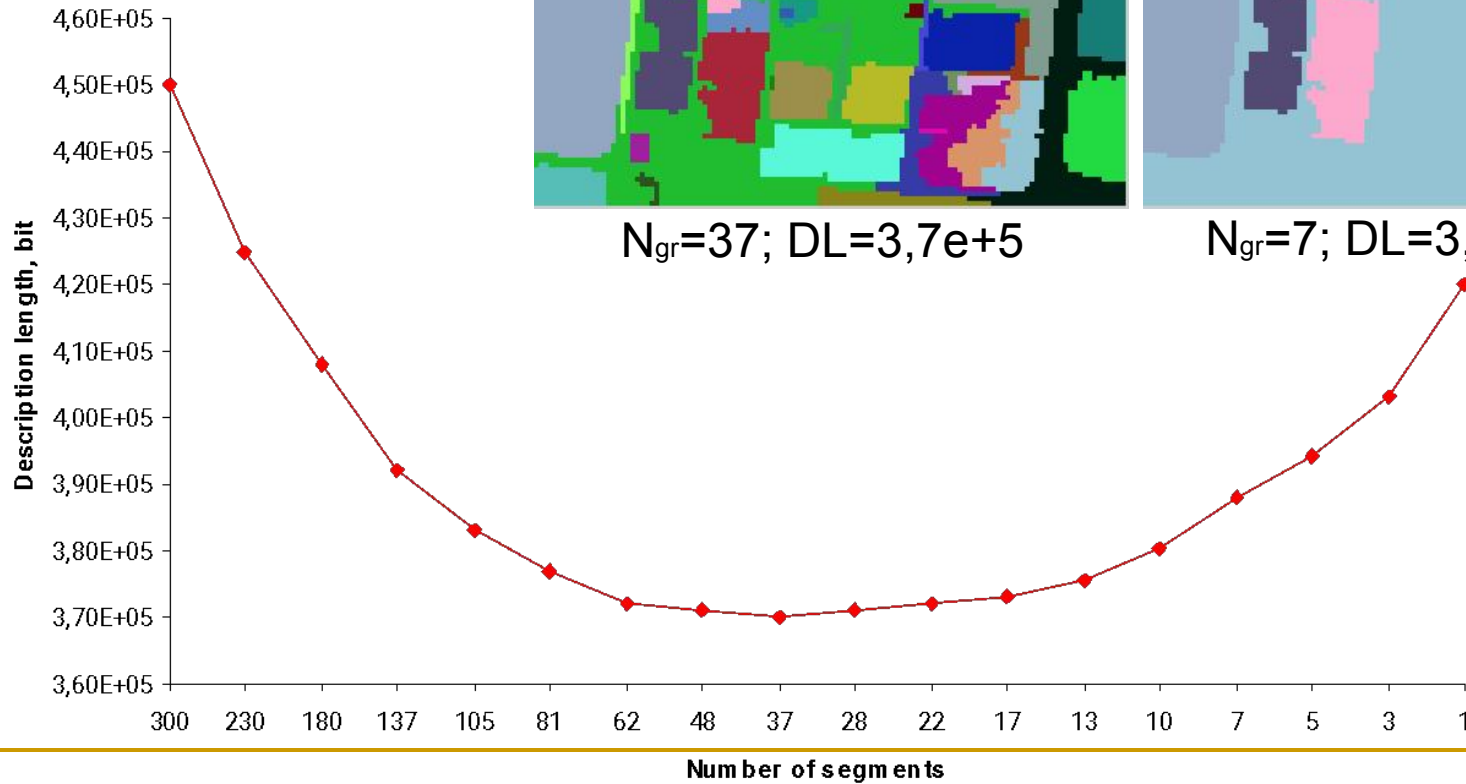
$N_{gr}=100$; $DL=3,8e+5$



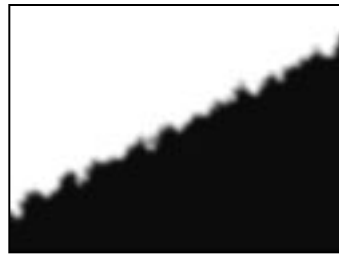
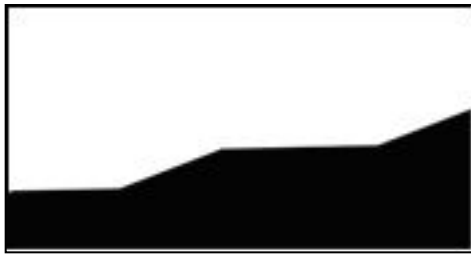
$N_{gr}=37$; $DL=3,7e+5$



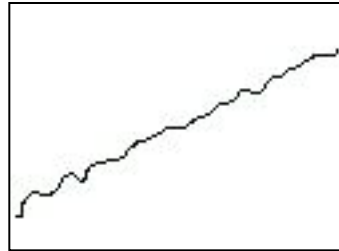
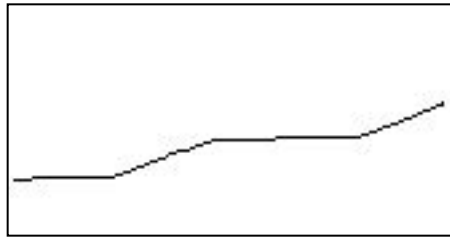
$N_{gr}=7$; $DL=3,9e+5$



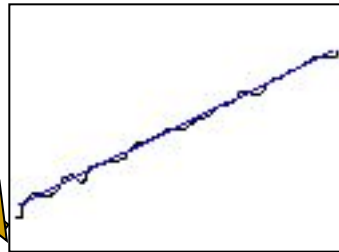
Contour segmentation



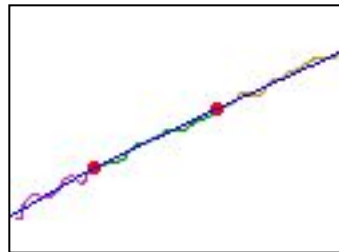
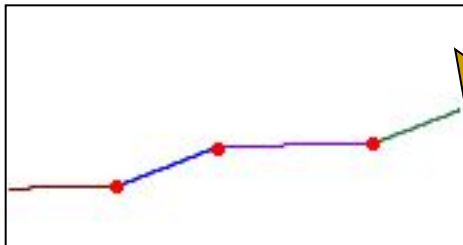
Images



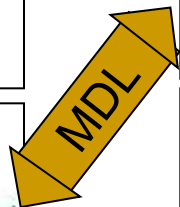
Extracted contours



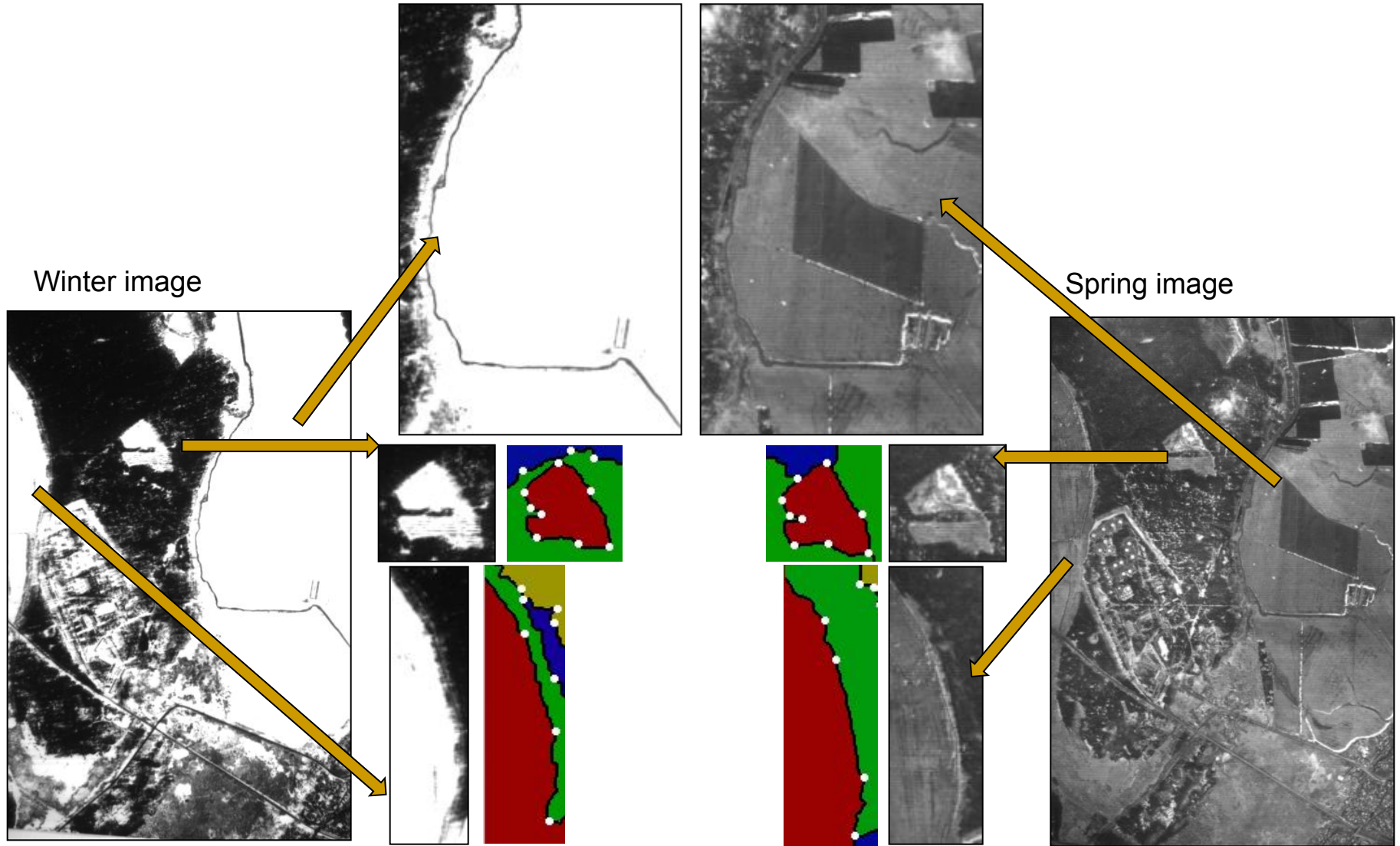
MSE-approximation with high threshold on dispersion



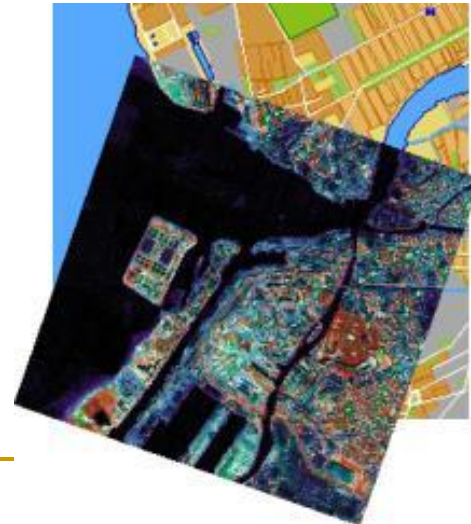
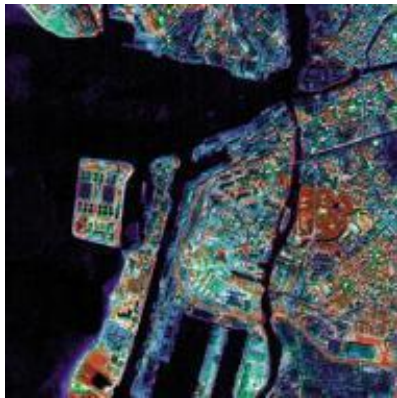
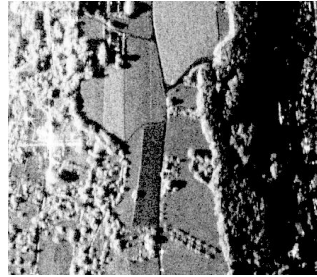
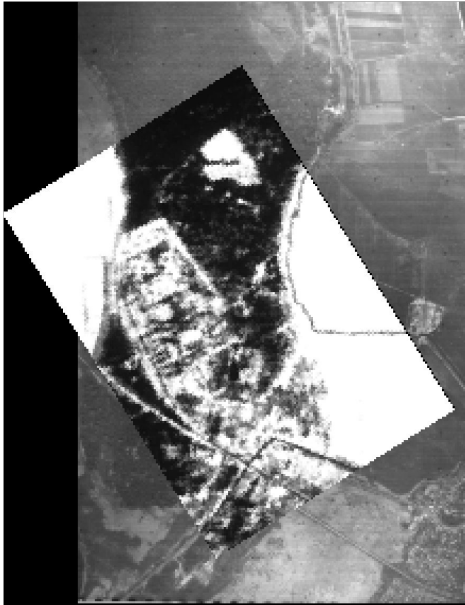
MSE-approximation with low threshold on dispersion



Full solution of invariant image matching



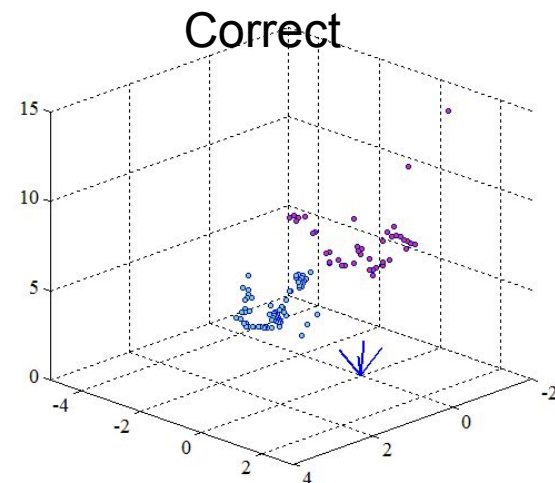
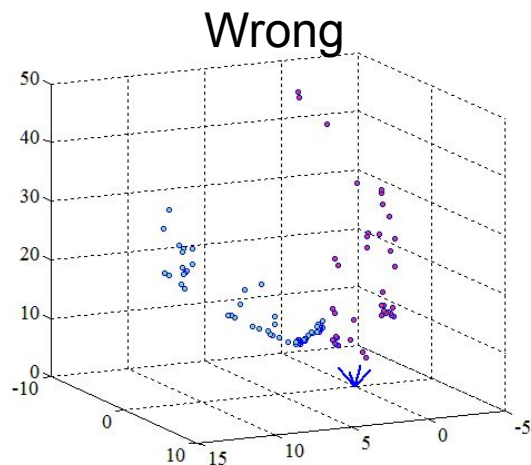
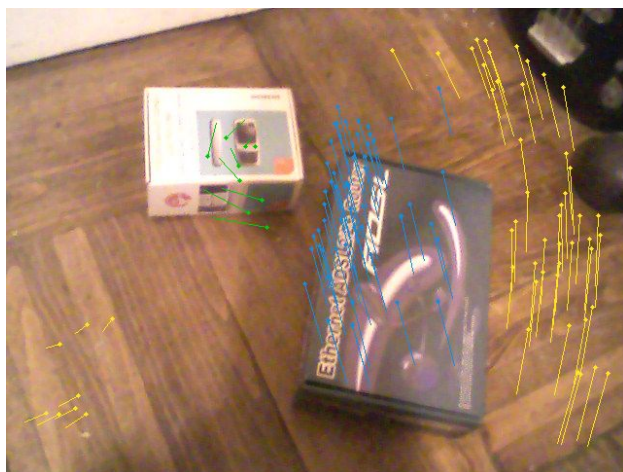
Successful matching



More applications of MDL



Correct separation into clusters for keypoint matching in dynamic scenes



Essential for correct estimation of a dynamic scene structure

A.N. Averkin, I.P. Gurov, M.V. Peterson, A.S. Potapov. Spectral-Differential Feature Matching and Clustering for Multi-body Motion Estimation // Proc. MVA2011 IAPR Conference on Machine Vision Applications. 2011. June 13-15, Nara, Japan. P. 173–176.

Various applications of MDL

- Pattern recognition, etc.:
 - Support-vector machines;
 - Discrimination functions;
 - Gaussian mixtures;
 - Decision forests;
 - ICA (as a particular case of MDL)
 - ...
- Image analysis
 - Segmentation;
 - Object recognition and image matching;
 - Optical flow estimation;
 - Structural description of images;
 - Changes detection;
 - ...
- Learning in symbolic domains, etc.

But wait... what about theory?

- MDL principle is used loosely
- Description lengths are calculated within heuristically defined coding schemes
- Success of a method is highly determined by the utilized coding scheme
- Is there some theory that overcomes this arbitrariness?

The theory behind MDL

- Algorithmic information theory

$$K_U(D) \square \min_H [I(H) \mid U(H) \square D], H^* \square \operatorname{argmin}_H [I(H) \mid U(H) \square D]$$

- U – universal Turing machine
- K – Kolmogorov complexity,
- $I(H)$ – length of program H
- H^* – best description/model of data D

- Two-part coding:

$$H^* \square \operatorname{argmin}_H [I(H) \square K(D \mid H)] \square \text{ if full model is separated into two parts}$$

OR

$$H^* \square \operatorname{argmin}_H [I(H) \square \log P(U(H) \square D)] \square \text{ if } H \text{ is probabilistic program}$$

- UTM defines the universal model space

Universal prediction

- Solomonoff's algorithmic probabilities
 - Prior probability

$$P_U(\alpha) = \sum_{p:U(p)=\alpha} 2^{-l(p)}$$

- Predictive probability

$$P_U(\alpha' | \alpha) = P_U(\alpha\alpha') / P_U(\alpha)$$

- Universal distribution of prior probabilities dominates (with multiplicative factor) over any other distribution
- Bayesian prediction with the use of these priors converges in limit with prediction based on usage of true distribution

Universality of the algorithmic space

3.1415926535 8979323846 2643383279 5028841971 6939937510 5820974944 5923078164 0628620899
8628034825 3421170679 8214808651 3282306647 0938446095 5058223172 5359408128 4811174502
8410270193 8521105559 6446229489 5493038196 4428810975 6659334461 2847564823 3786783165
2712019091 4564856692 3460348610 4543266482 1339360726 0249141273 7245870066 0631558817
4881520920 9628292540 9171536436 7892590360 0113305305 4882046652 1384146951 9415116094
3305727036 5759591953 0921861173 8193261179 3105118548 0744623799 6274956735 1885752724
8912279381 8301194912 9833673362 4406566430 8602139494 6395224737 1907021798 6094370277
0539217176 2931767523 8467481846 7669405132 0005681271 4526356082 7785771342 7577896091
7363717872 1468440901 2249534301 4654958537 1050792279 6892589235 4201995611 2129021960
8640344181 5981362977 4771309960 5187072113 4999999

```
int a=10000,b,c=8400,d,e,f[8401],g;
main() {for(;b-c;)f[b++]=a/5;
for(;d=0,g=c*2;c-=14, printf("%.4d",e+d/a),e=d%a)
for(b=c;d+=f[b]*a,f[b]=d%--g,d/=g--,--b;d*=b);}
```

By D.T. Winter

Grue Emerald Paradox

- Hypothesis No. 1: all emeralds are green
- Hypothesis No. 2: all emeralds are greu
(that is green before 2050, and blue after this time)

- Likelihood of observation data equals
- How can we calculate prior probabilities of these two hypotheses?

Is it possible to ground prior probabilities?

- Probability theory allows to deduce one probability from another. But what are the initial probabilities?
- Universal priors work

Methodological usefulness

- Theory of universal induction answers the questions
 - What is the source of overlearning/ overfitting/ oversegmentation, etc.
 - Why is any new narrow learning method “yet another classifier”
 - Why are feed forwards neural networks not really “universal approximators”
 - And at the same time, why is “no free lunch theorem” not true

Gap between universal and pragmatic methods

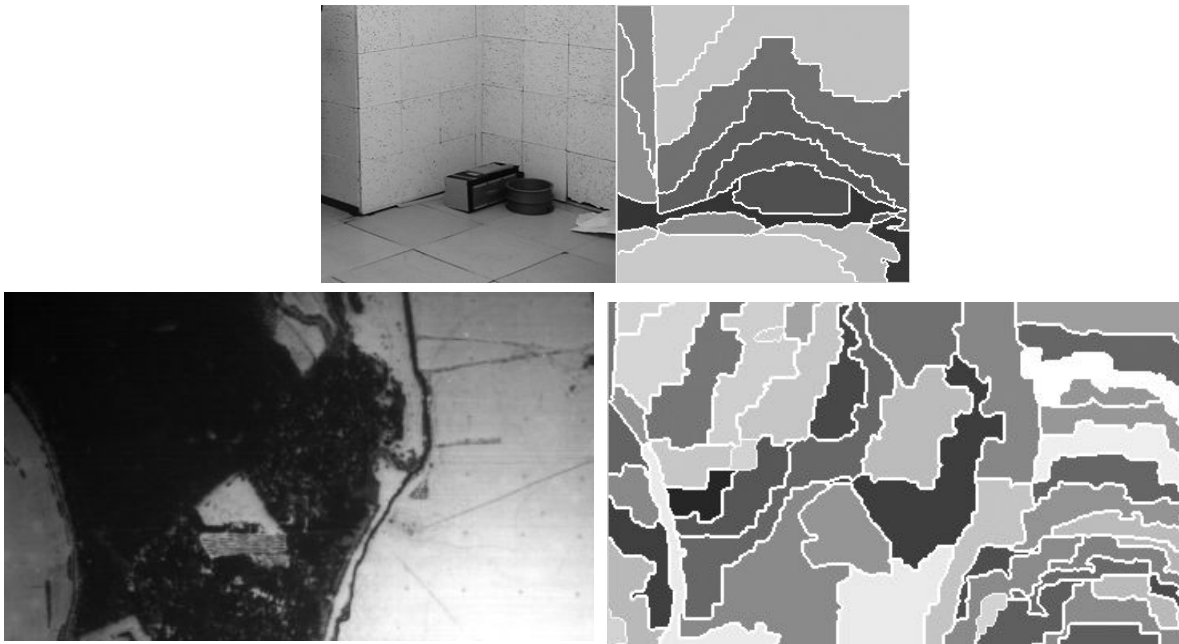
- Universal methods
 - can work in arbitrary computable environment
 - incomputable or computationally infeasible
 - approximations are either inefficient or not universal
 - Practical methods
 - work in non-toy environments
 - set of environments is highly restricted
- ⇒ Bridging this gap is necessary

Choice of the reference UTM

- Unbiased AGI cannot be practical and efficient
- Dependence of the algorithmic probabilities on the choice of UTM appears to be very useful in order to put any prior information and to reduce necessary amount of training data
- UTM contains prior information
 - ⇒ UTM can be optimized to account for posterior information

Limitations of narrow methods

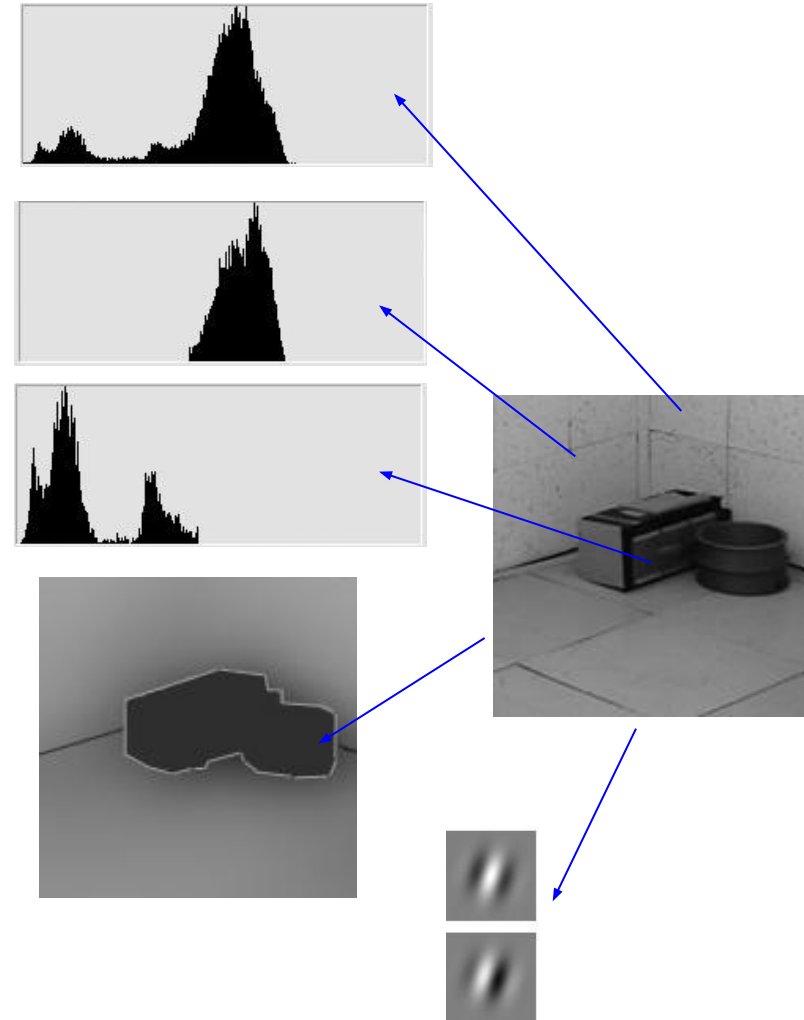
- Brightness segmentation can fail even with the MDL criterion



Essentially incorrect segments

More complex models...

1. Image is described as a set of independent and identically distributed samples of random variable (no segmentation).
2. Image is divided into regions; brightness values described independently within each region.
3. Second order functions are fit in each region, and brightness residuals are described as iid random variables.
4. Mixes of Gabor functions are used as regression models.

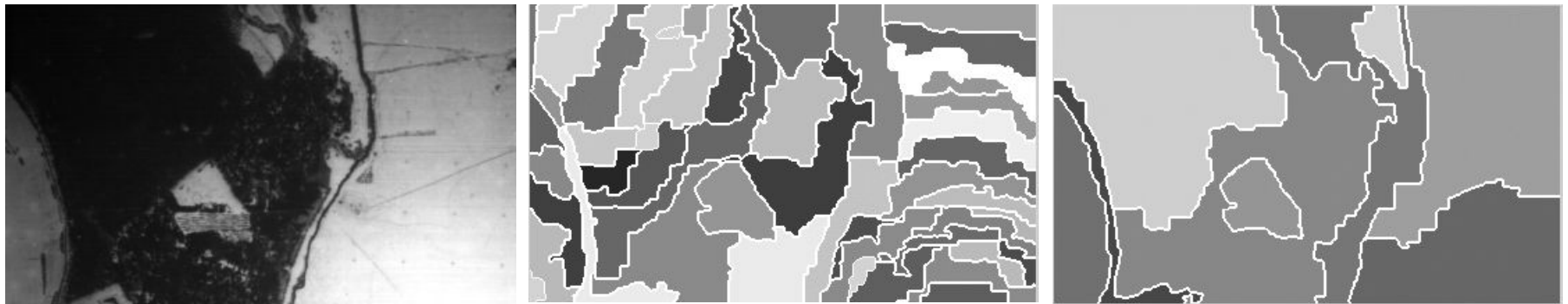
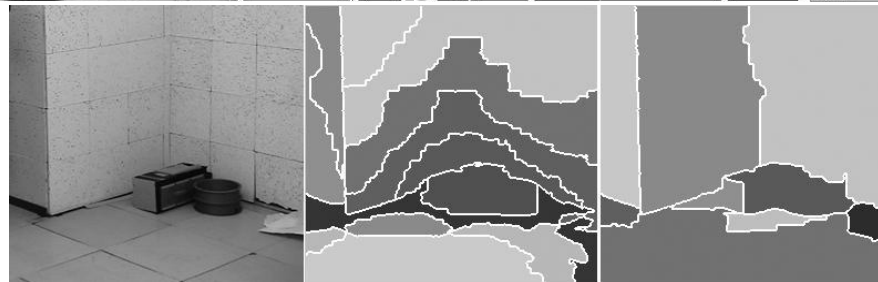
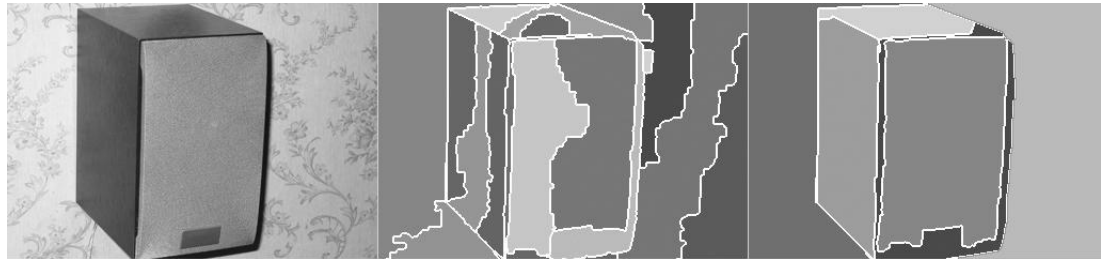


Comparison

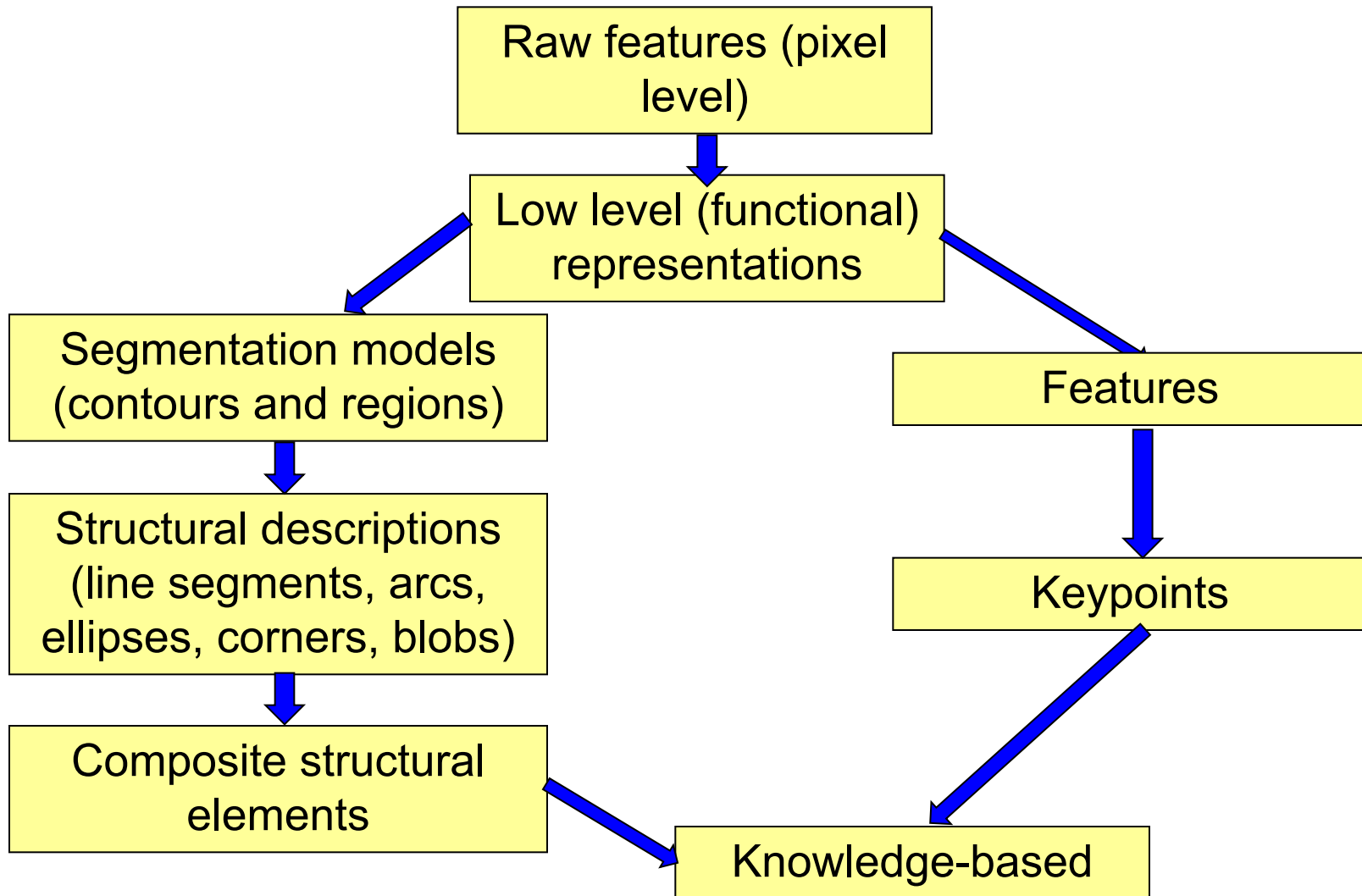
Images

Brightness
entropy

Regression
models

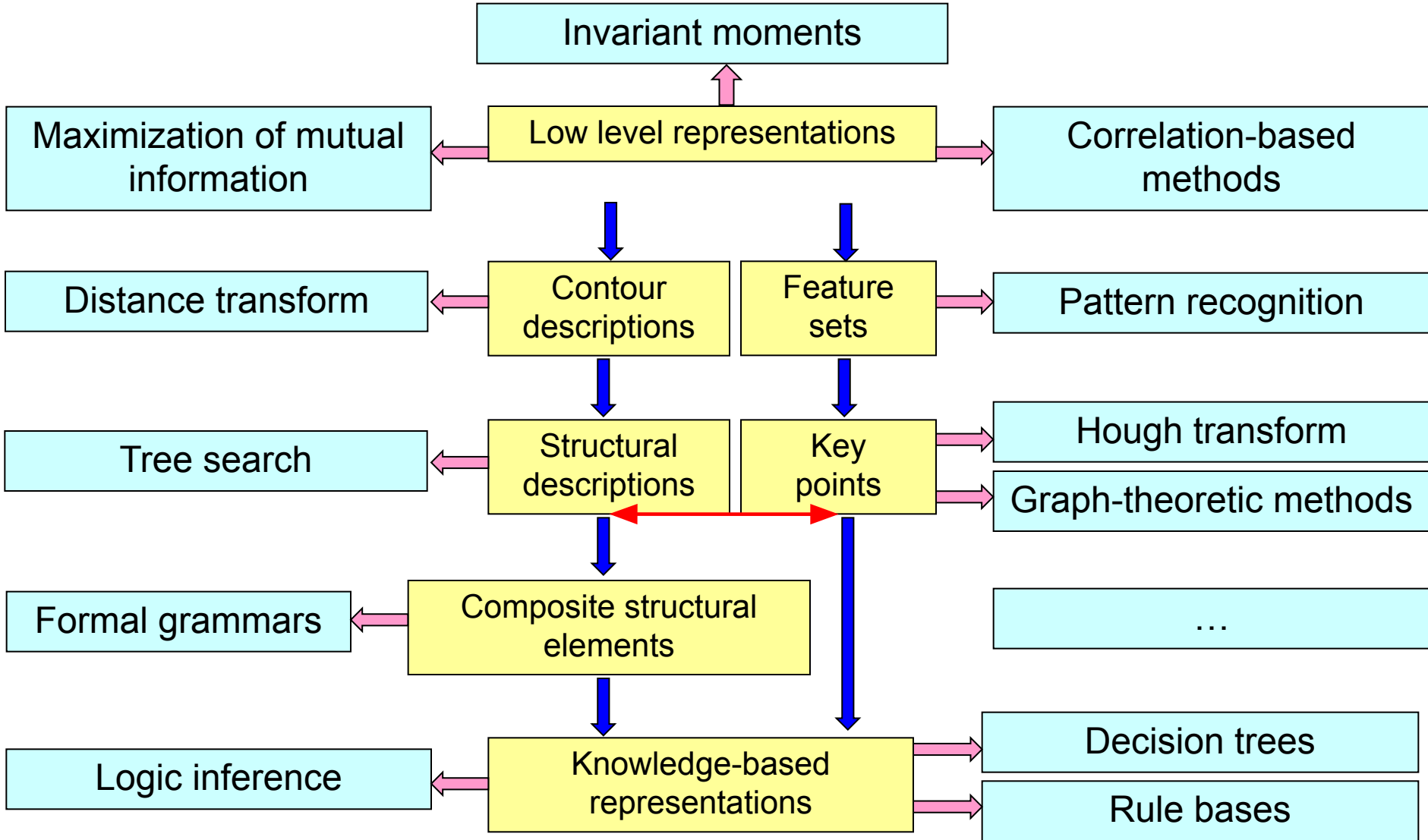


Classes of image representations*



*Marr, D.: Vision: A Computational Investigation into the Human Representation and Processing of Visual Information. MIT Press (1982).

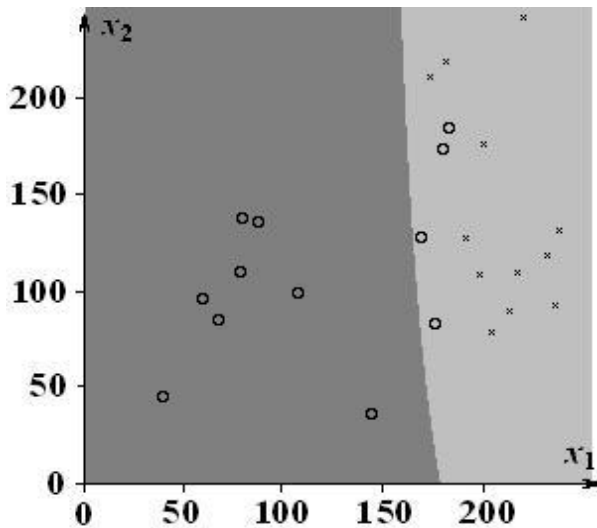
Example: image matching



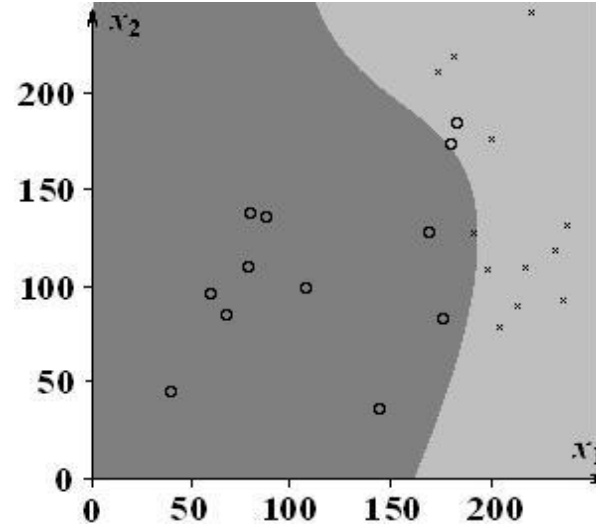
But again... what about theory?

- MDL principle is used loosely
 - Description lengths are calculated within heuristically defined coding schemes
 - Success of a method is highly determined by the utilized coding scheme
 - In computer vision and machine learning, some representation is used in every method
- ⇒ **But how to construct the best representation?**
- ⇒ Representations correspond to 'coding schemes' in MDL applications. They should also be constructed on the base of strict criterion
- ⇒ **But from what space and how?**

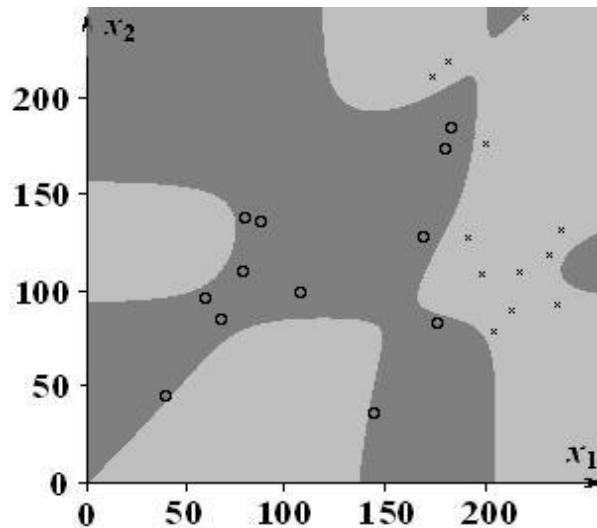
Polynomial decision function



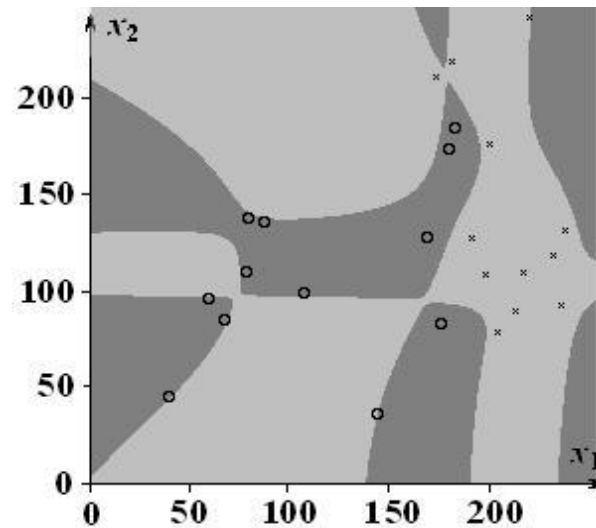
%(learn)=11.1
%(test)=5.4
L = 31.2 bit
Np=4



%(learn)=2.8
%(test)=3.6
L = 30.9 bit
Np=9



%(learn)=0.0
%(test)=8.6
L = 41.4 bit
Np=16

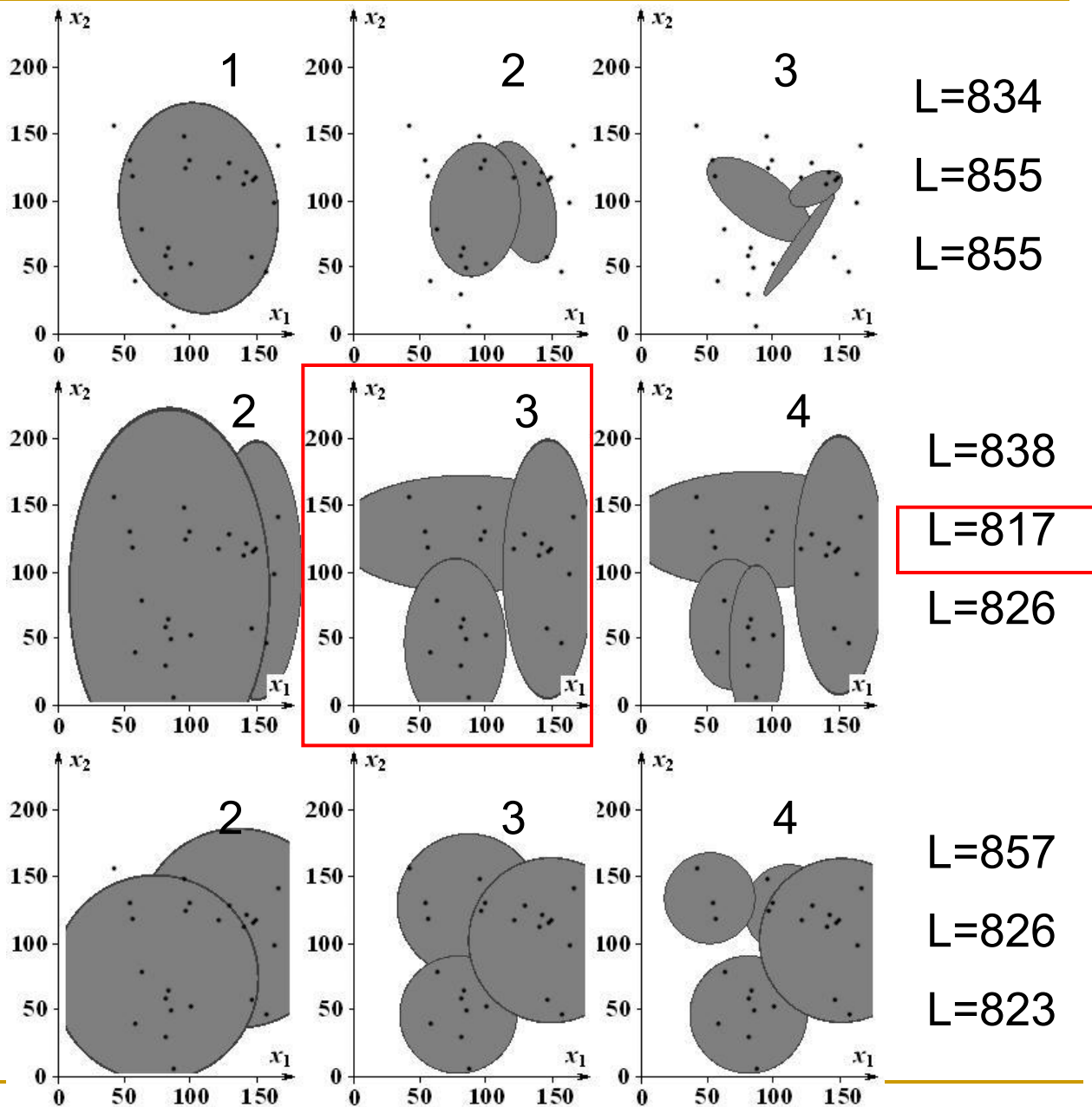


No outliers

Worst generalization!

%(learn)=0.0
%(test)=18.4
L = 62.0 bit
Np=25

Choosing between mixtures with different number of components and restrictions laid on the covariance matrix of normal distribution



Again, heuristic coding schemes

- Let's switch back to theory

Universal Mass Induction

- Let $\{x_i\}_{i=1}^n$ be the set of strings
- An universal method cannot be applied to mass problems since typically

$$K_U(x_1x_2\dots x_n) \ll \sum_{i=1}^n K_U(x_i)$$

where K is Kolmogorov complexity on universal machine U

- However, $K_U(x_1x_2\dots x_n) \approx \min_S \left(I(S) + \sum_{i=1}^n K_U(x_i | S) \right)$ can hold
- One can search for models $y_i^* = \operatorname{argmin}_{y: S(y)=x_i} I(y)$ for each x_i independently

within some best representation $S^* = \operatorname{argmin}_S \left(I(S) + \sum_{i=1}^n I(y_i^*) \right)$

Representational MDL principle

For example, image analysis tasks are mass problems: the same algorithm is applied to different images (or patterns) independently.

- Definition

Let *representation* for the set of data entities be such the program S for UTM U that for any data entity D the description H exists that $U(SH)=D$.

- Representational MDL principle

- The best image description has minimum length within given representation
- The best image representation minimizes summed description length of images from the given training set (and the length of representation itself).

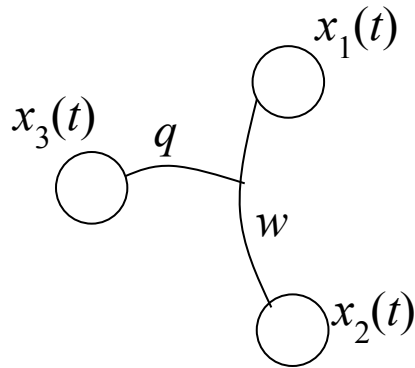
Main advantage: applicable to any type of representation; representation is included into general criterion as a parameter.

Possible usage of RMDL

- Synthetic pattern recognition methods*:
 - Automatic selection among different pattern recognition methods
- Selecting a representation that better fits the training sample from a specific domain either from a family of representations or from a fixed set of hand-crafted representations
- Improve data analysis methods for specific representations

RMDL for optimizing ANN formalisms

- Considered extension of ANN representation

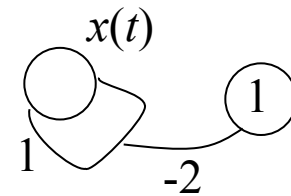


$$x'_2(t) = wx_1^{qx_3(t)+1}(t)$$

$$q = 0 \vee x_3 = 0 \Rightarrow x'_2(t) = wx_1(t)$$

$$qx_3 = 1 \Rightarrow x'_2(t) = wx_1^2(t)$$

$$qx_3 = -2 \Rightarrow x'_2(t) = wx_1^{-1}(t)$$

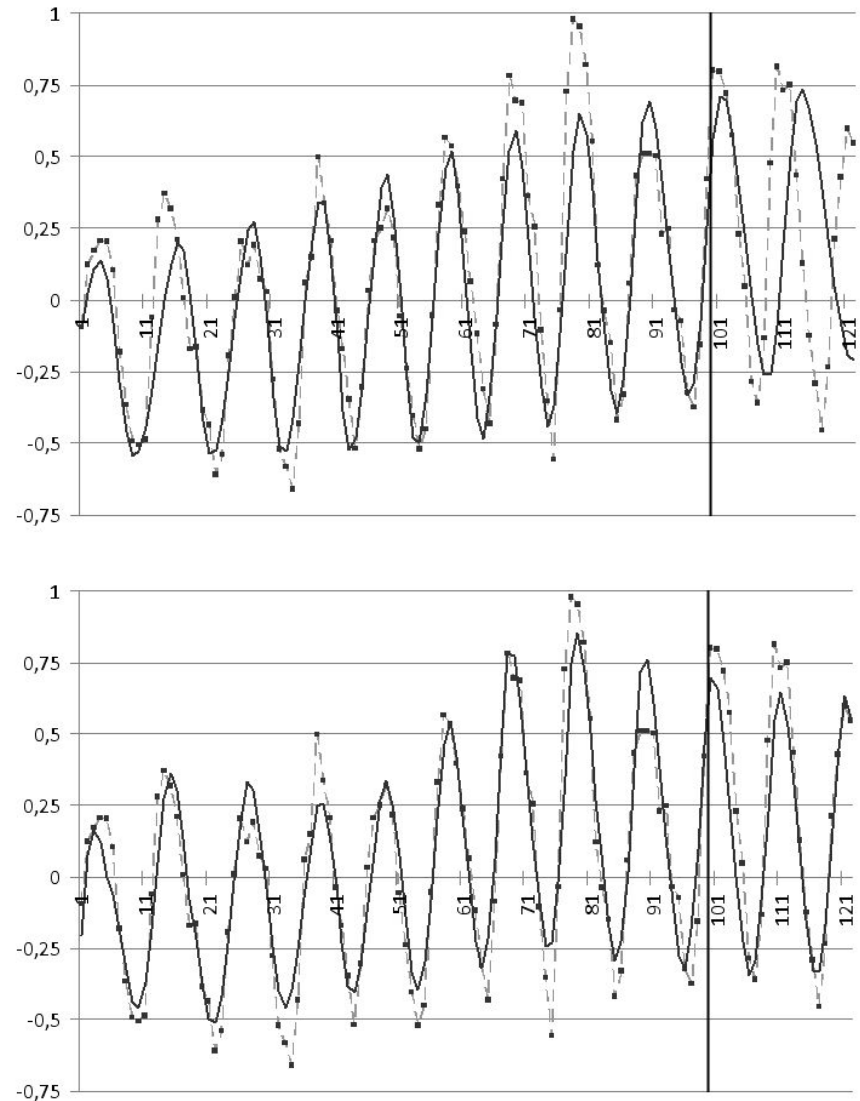


$$x'(t) = 1/x(t)$$

$$x(t) = \ln(t)$$

RMDL for optimizing ANN formalisms

- Experiments: Wolf annual sunspot time series
- Precision of forecasting depends on type of nonlinearity
- ANN with 4 neurons, 11 connections, and 2 second-order connections: MSE=220 (typical MSE: 214–625*)

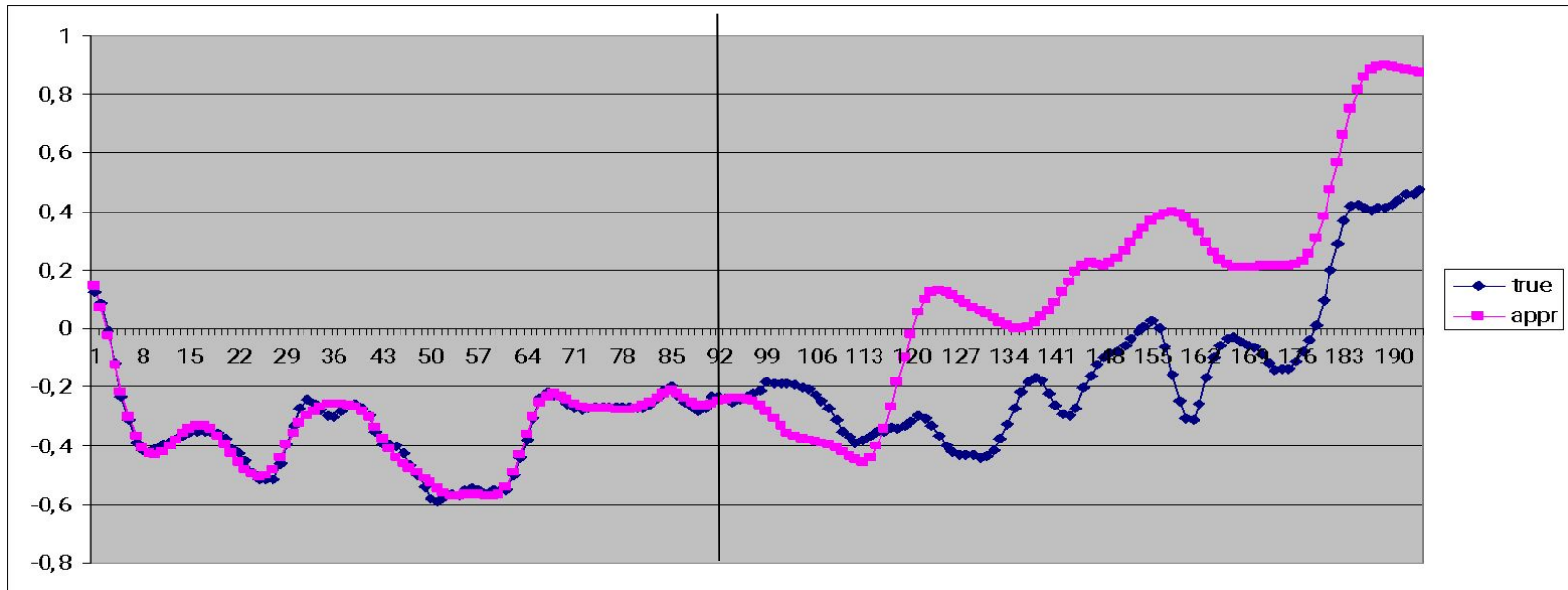


RMDL for optimizing ANN formalisms

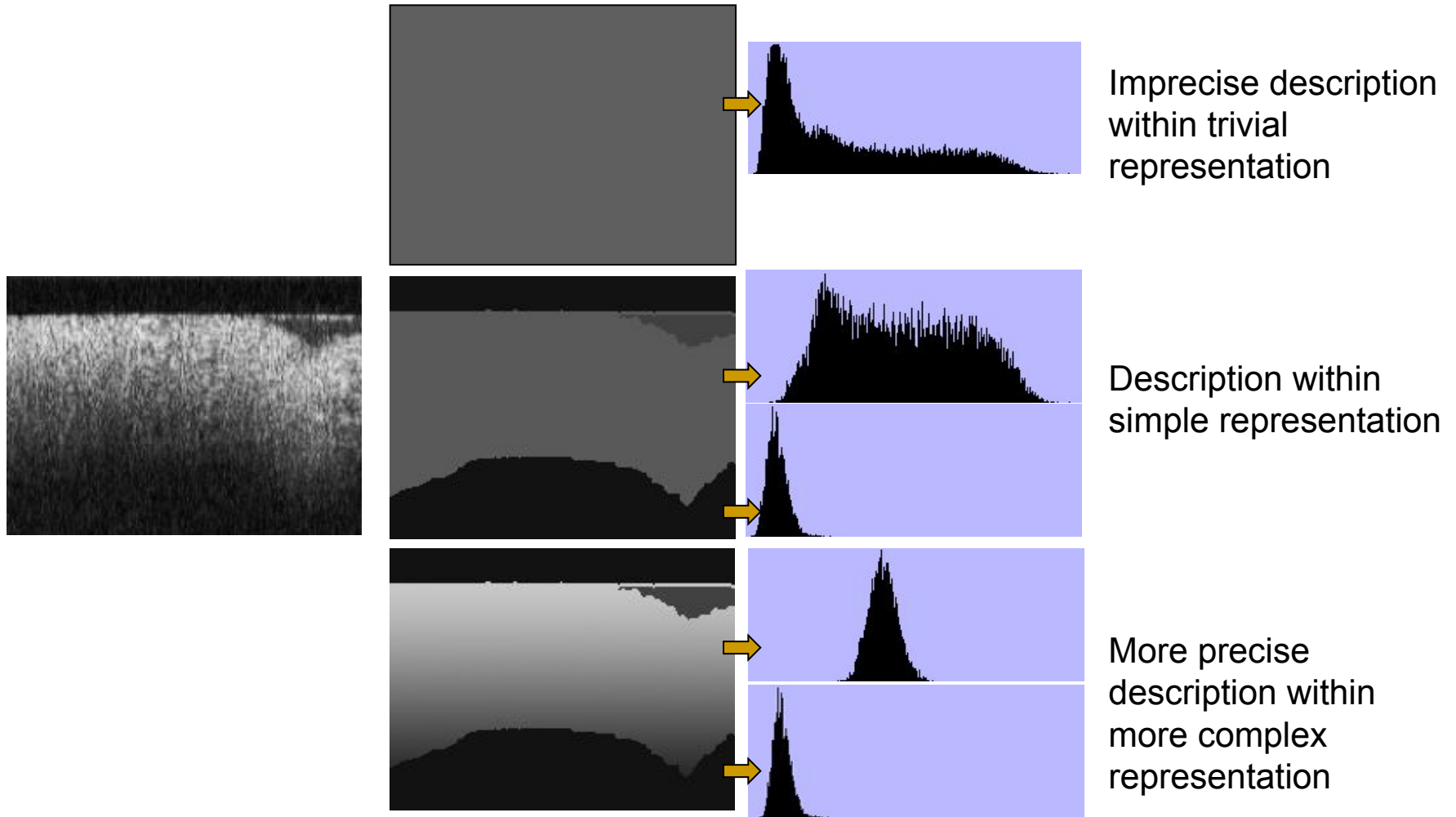
ANN type	RMDL, bits	error, %
Linear	651	15,8
Activation function	617	10,1
2 nd -order connections	608	9,9

Test: Financial time series

Although we obtained an agreement between the short-term prediction precision and the RMDL criterion in average, one can agree with the statement: “MSE and NMSE are not very good measures of how well the model captures the dynamics”



OCT image segmentation

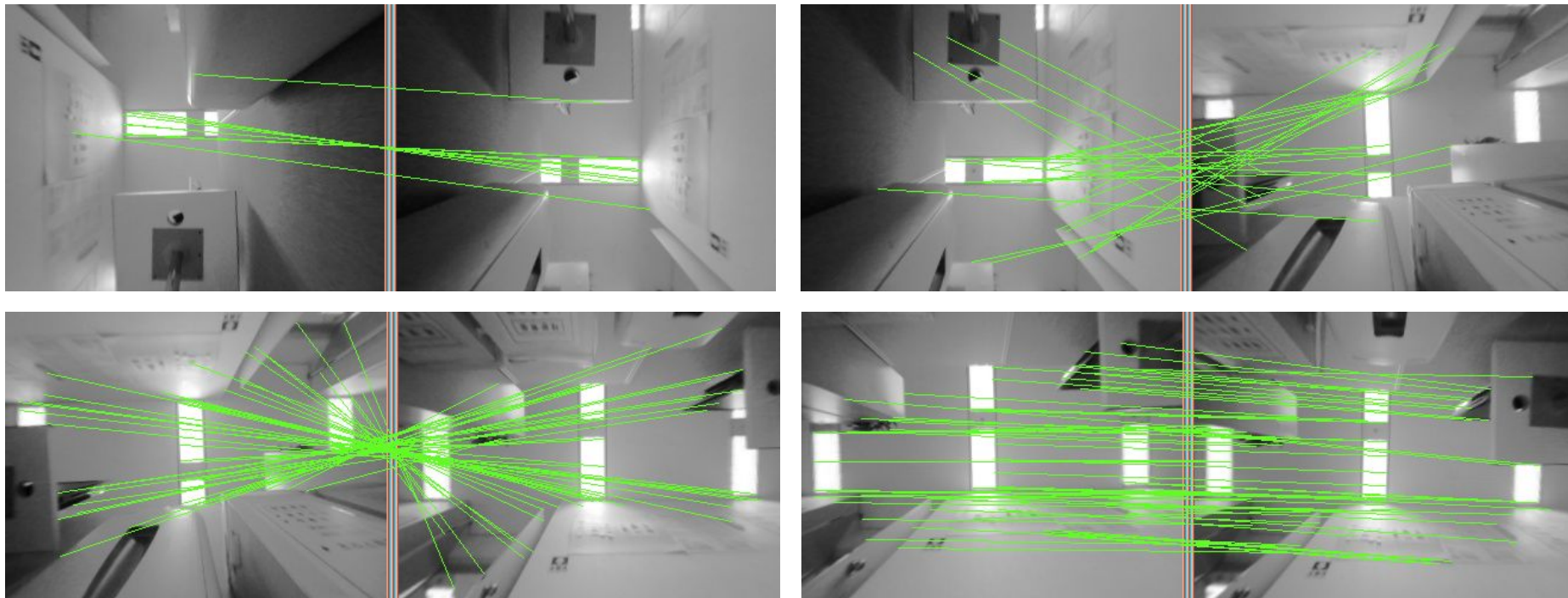


Segmentation results

			Description length, bits S-0: 212204 S-1: 184672 S-2: 175096	S-1: oversegmentation S-2: correct detection of layers
			Description length, bits S-0: 231201 S-1: 212268 S-2: 207864	S-1 and S-2 are almost the same (and plausible) detection of thin layers
			Description length, bits S-0: 235566 S-1: 219641 S-2: 215066	Differing segmentation results for a single thick layer (light absorption with depth causes regular reduction of brightness). Some inclusions are not detected.
			Description length, bits S-0: 236421 S-1: 213015 S-2: 206204	S-1: odd layer is detected and inclusion is missed S-2: plausible results of segmentation

Application to image feature learning

Training set with preliminarily matched key points using predefined hand-crafted feature transform



Example of some found linear feature transforms

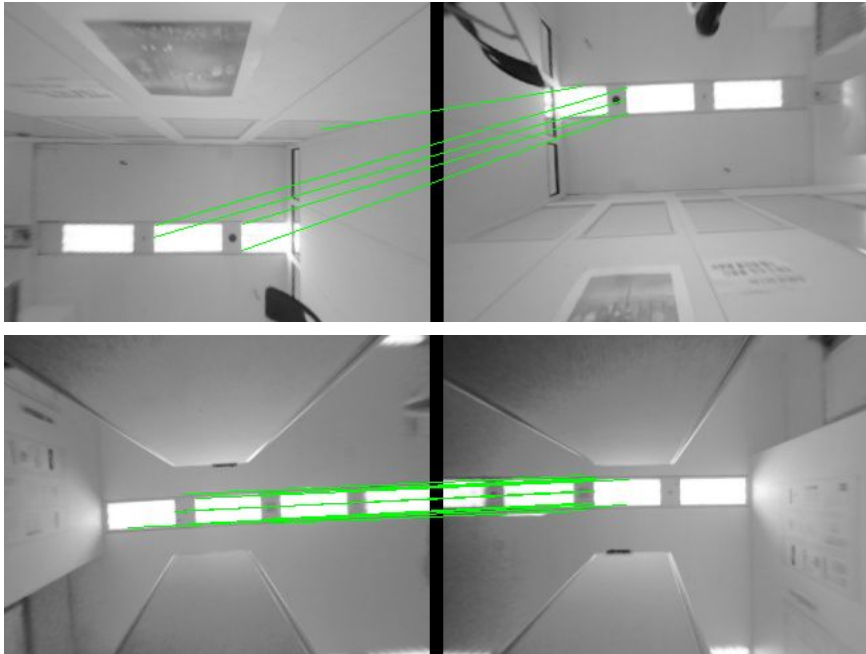


Example of some feature transforms for another environment

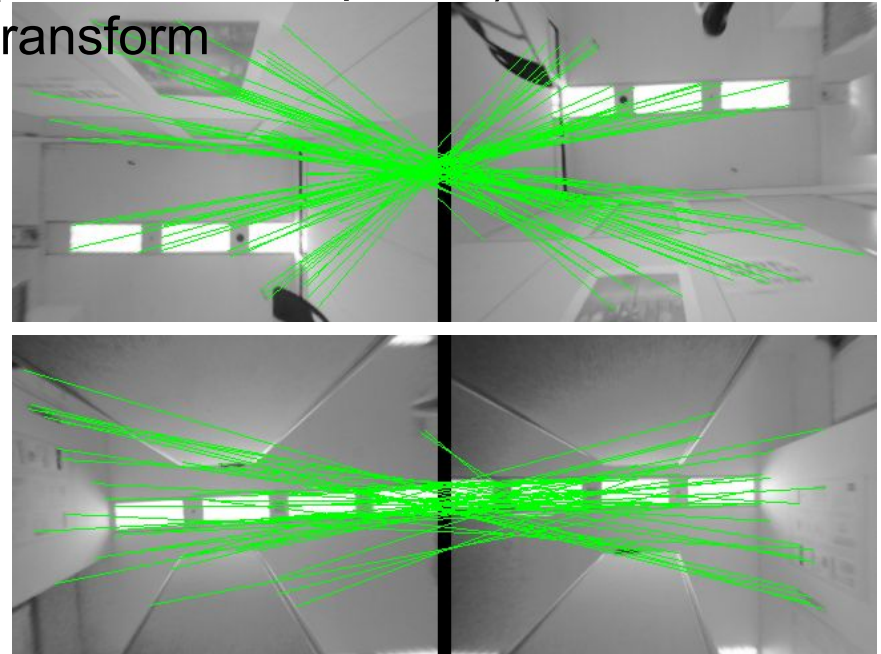


Results

Matching with predefined
hand-crafted feature transform

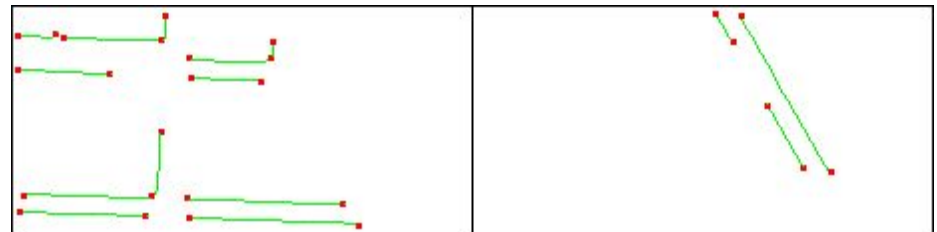
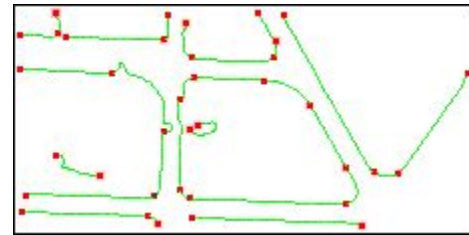
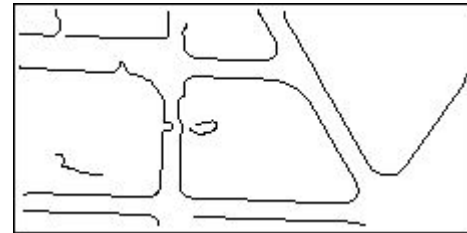
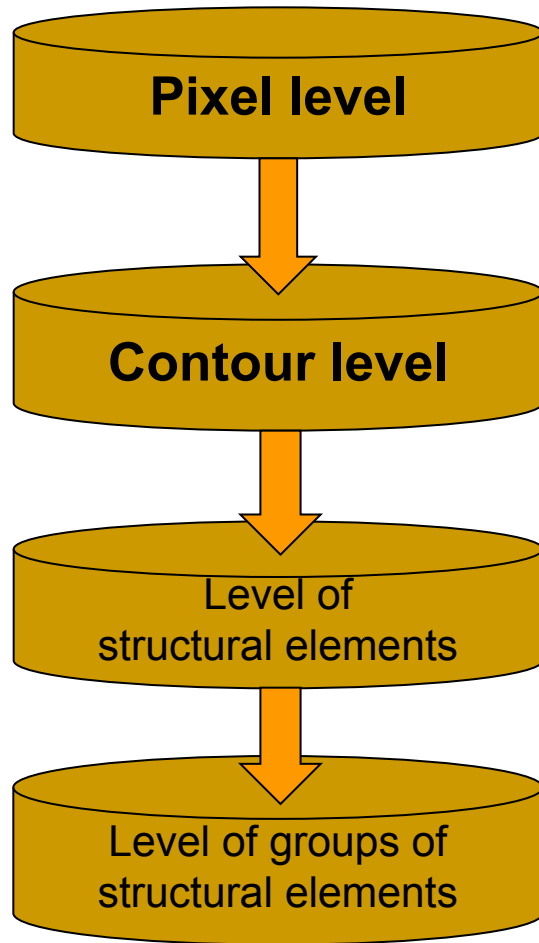


Matching with learned
(environment-specific) feature
transform

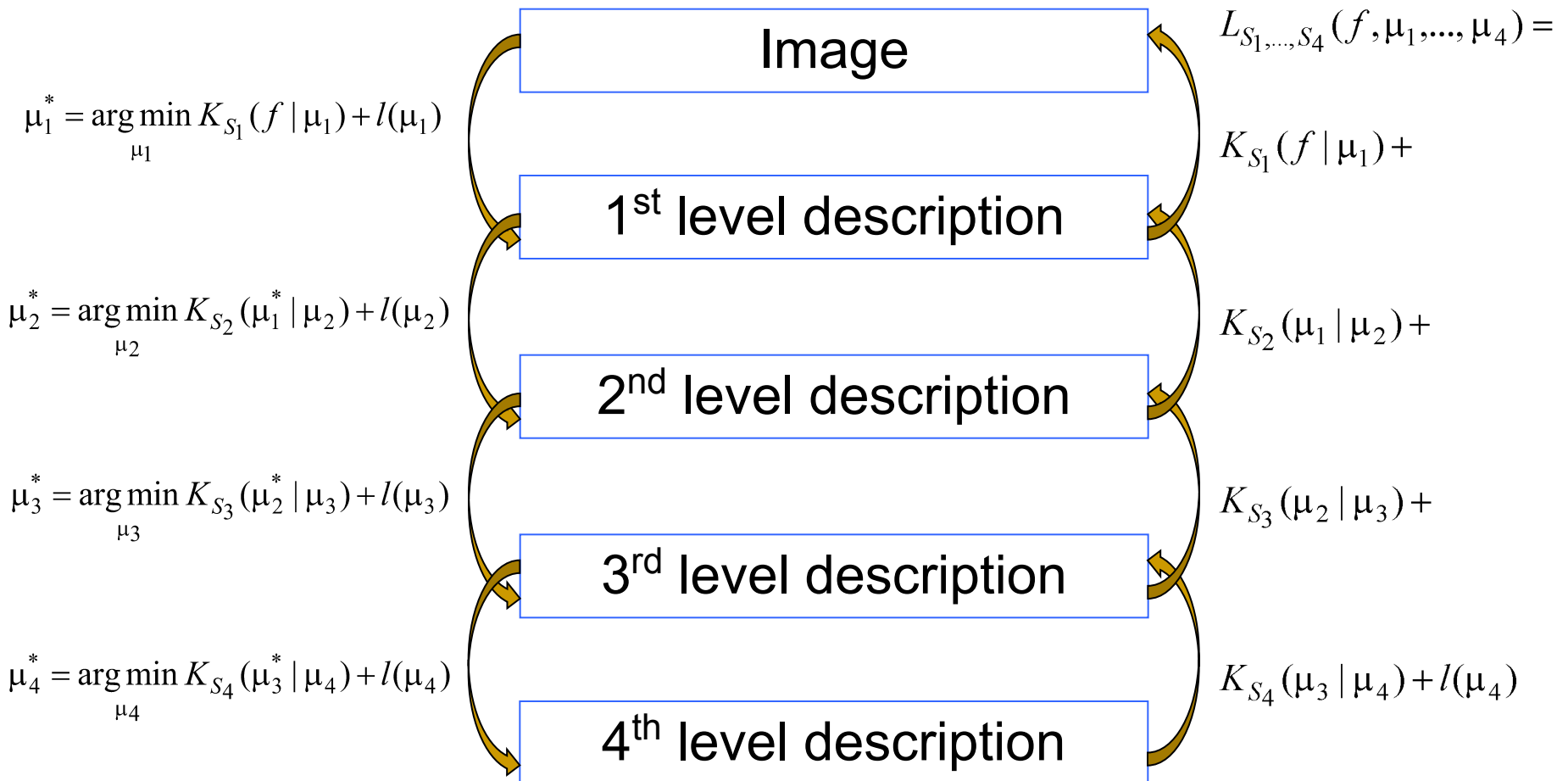


~50% of failures with predefined features were matched successfully
with learned features (new images of the same environment were used)

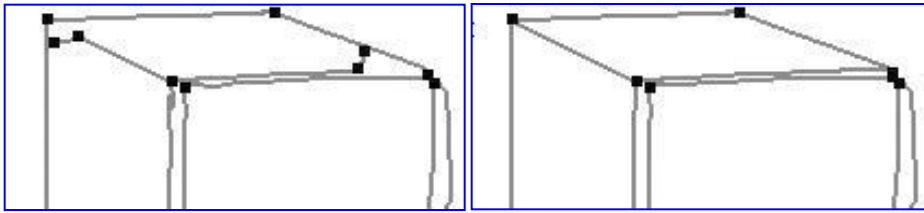
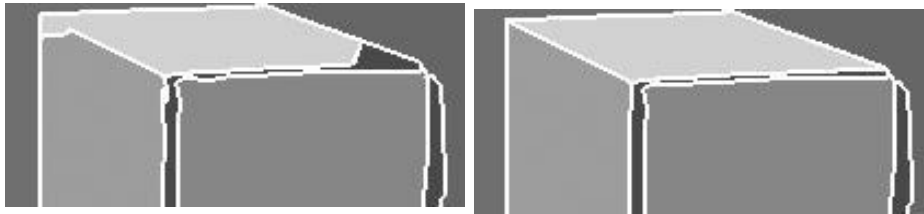
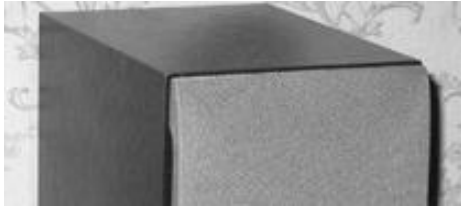
Analysis of hierarchical representations



Adaptive resonance



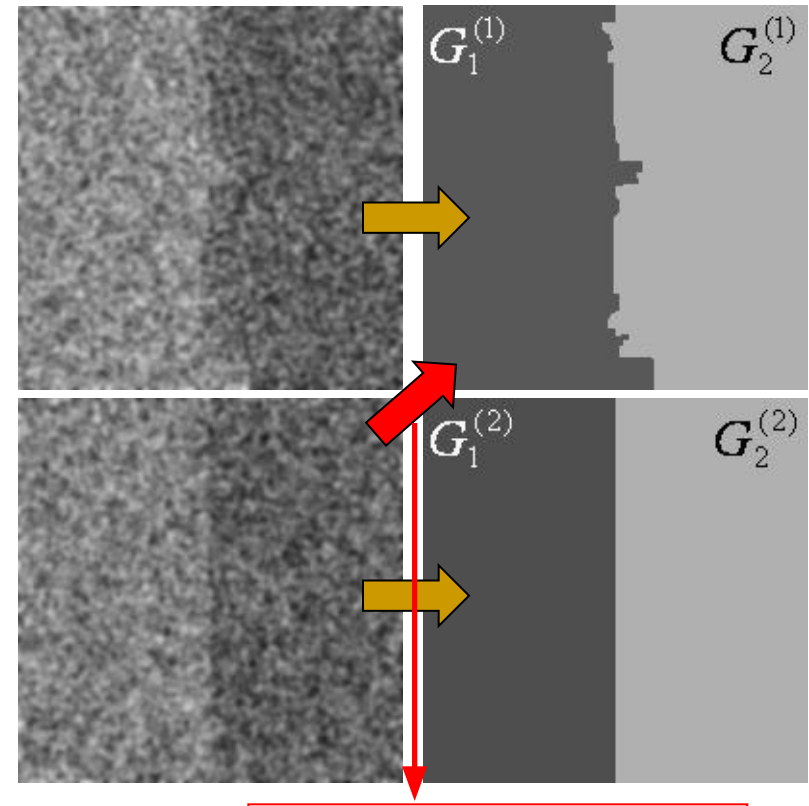
Implications



Independent
optimization of
descriptions

Usage of integral
description length

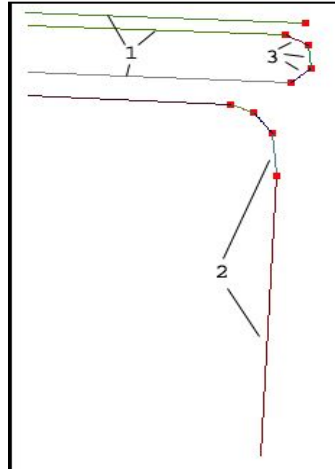
$$L_{S_1^{(H)}}(f, \mu_1, \mu_2, \mu_3, \mu_4) = K_{S_1^{(1)}}(f | \mu_1) + K_{S_1^{(1)}}(\mu_1 | \mu_2) + \\ + K_{S_2^{(2)}}(\mu_2 | \mu_3) + K_{S_1^{(3)}}(\mu_3 | \mu_4) + K_{S_0^{(4)}}(\mu_4).$$



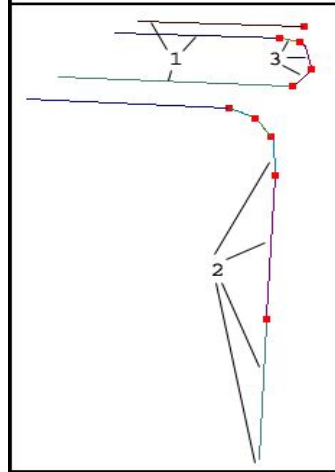
Without resonance

Adaptive resonance: matching as construction of common description

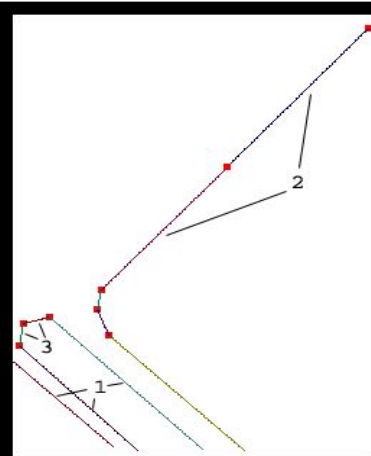
Initial structural elements of the first image



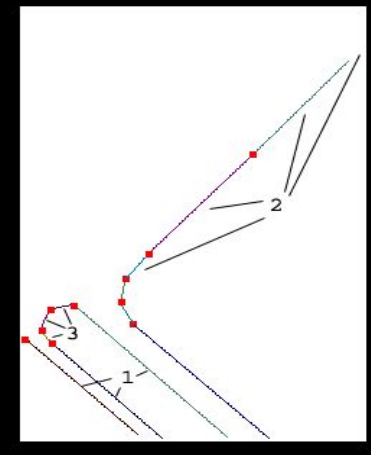
Fixed structural descriptions: same for both images



Initial structural elements of the second image



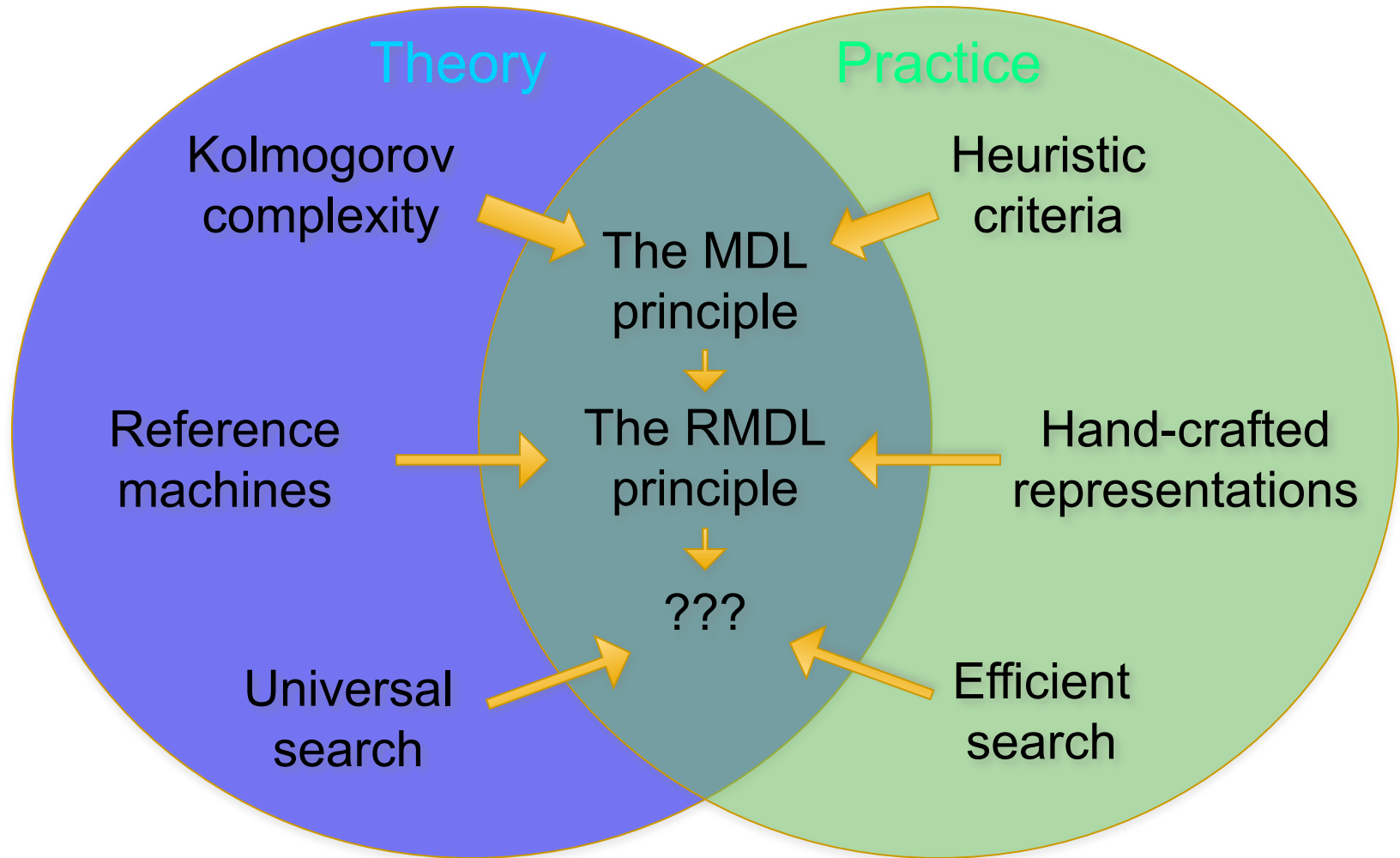
These descriptions slightly less precise, but w.r.t. images, but only one of them can be used instead of two



Learning representations

- Very difficult problem in Turing-complete settings
- Successful methods use efficient search and restricted families of representations
- Deep learning
 - Not universal
 - Compact (one-level ANNs should be exponentially larger than multi-level ANNs to represent some concepts => particular case of RMDL)
 - Higher expressive power or more efficient search than those of former methods

What is still missing?



Key Idea

- Humans create narrow methods, which efficiently solve arbitrary recurring problems
- Generality should be achieved not by a single uniform method solving any problem in the same fashion, but by automatic construction of (non-universal) efficient methods
- Program specialization is the appropriate concept*, which relates general and narrow intelligence methods
- However, no analysis of possible specialization of concrete models of universal intelligence has been given yet.

Program Specialization

- Let $p_L(x,y)$ be some program (in some language L) with two arguments
- Specializer $spec_R$ is such program (in some language R) accepting p_L and x_0 that

$$(\forall y)spec_R(p_L, x_0)(y) = p_L(x_0, y)$$

- $spec_R(p_L, x_0)$ is the result of deep transformation of p_L that can be much more efficient than $p(x_0, \cdot)$

Futamura-Turchin projections

$$(\forall x)spec_R(intL, p_L)(x) = intL(p_L, x)$$

$$(\forall p_L, x)spec_R(spec_R, intL)(p_L)(x) = intL(p_L, x)$$

$$(\forall intL)spec_R(spec_R, spec_R)(intL) = comp_{L \rightarrow R}$$

Specialization of Universal Induction

- Universal mass induction consists of two procedures
 - Search for models

$$MSearch(S, x_i) \rightarrow y_i^* = \operatorname{argmin}_{y: S(y)=x_i} I(y)$$

- Search for representations

$$RSearch(x_1, \dots, x_n) \rightarrow S^* = \operatorname{argmin}_S \left(I(S) + \sum_{i=1}^n I(y_i^*) \right)$$

- $MSearch(S, x)$ is executed for different x with same S
- This search cannot be non-exhaustive for any S , but it can be efficient for some of them
- One can consider computationally efficient projection
 $spec(MSearch, S): (\forall x) spec(MSearch, S)(x) = MSearch(S, x)$

Approach to Specialization

- Direct specialization of $MSearch(S, x)$ w.r.t. some given S^*
 - No general techniques for exponential speedup exists
 - And how to get S' ? $RSearch$ is still needed
- Find $S' = spec(MSearch(S, x), S^*)$ simultaneously with S^*

Main properties of S, S' : $(\forall x)S(S'(x)) = x$

$$I(S) + \sum_i I(S'(x_i)) \rightarrow \min$$

- S is a generative representation (decoding)
- S' is a descriptive representation (encoding)
 - S' is also the result of specialization of the search for generative models, so in general it can include some sort of optimized search
- Simultaneous search for S and S' will be referred to as SS' -search

Conclusion

- Attempts to build more powerful practical methods led us to utilization of the MDL principle that was heuristically applied for solving many tasks
- The MDL principle is a very useful tool for introducing model selection criteria free from overfitting in the tasks of image analysis and pattern recognition
- We introduced the representational MDL principle to bridge the gap between universal induction and practical methods and used it to extend practical methods
- The remaining difference between universal and practical methods is in search algorithms. Specialization of universal search is necessary to automatically produce efficient methods

Thank you for attention!

Contact: potapov@aideus.com