

Learning End-to-End

Thilo Stadelmann & Thierry Musy

Datalab Seminar, 24. June 2015, Winterthur



DISCLAIMER



Today's talk is about the basic ideas of a single, inspiring,
industry-proven paper from the nineties...



LeCun et al., “Gradient-Based Learning Applied to Document Recognition”, 1998

Agenda

Challenge

The bigger picture of ML – Sequential Supervised Learning as an example where simplistic paradigms fail

Proposed Solution

Global learning – Graph Transformer Networks – Example: Digit string recognition with heuristic oversegmentation – Advantage of discriminative training

Use Case

Error propagation through multiple modules – Check reading application



Conclusions

Summary – Outlook to other projects – Remarks

CHALLENGE: SEQUENTIAL SUPERVISED LEARNING

Machine Learning

Wikipedia on «Learning», 2015:

«...the act of **acquiring** new, or **modifying** and reinforcing, existing **knowledge**, behaviors, **skills**, **values**, or preferences and **may involve synthesizing** different types of information.»

A. Samuel, 1959:

«...gives computers the **ability** to learn **without being explicitly programmed.**»

“do something”

T.M. Mitchell, 1997:

«...if its **performance** at tasks in T, as measured by P, **improves with experience E.**»

→ In practice: Fitting parameters of a function to a set of data.
(data usually handcrafted, function chosen heuristically)



A landmark work in the direction of «learning»

LeCun et al., “Gradient-Based Learning Applied to Document Recognition”, 1998



Outline

- Gradient-Based ML ✓
- Convolutional Neural Networks ✓
- Comparison with other Methods
- Multi-Module Systems & Graph Transformer Networks
- Multiple Object Recognition & Heuristic Oversegmentation
- Space Displacement Neural Networks
- GTN's as General Transducers
- On-Line Handwriting Recognition System
- Check Reading System

Standard and Sequential Supervised Learning

Supervised Learning

Incanter Dataset				
Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
5.1	3.5	1.4	0.2	setosa
4.9	3.0	1.4	0.2	setosa
4.7	3.2	1.3	0.2	setosa
4.6	3.1	1.5	0.2	setosa
5.0	3.6	1.4	0.2	setosa
5.4	3.9	1.7	0.4	setosa
4.6	3.4	1.4	0.3	setosa
5.0	3.4	1.5	0.2	setosa
4.4	2.9	1.4	0.2	setosa
4.9	3.1	1.5	0.1	setosa
5.4	3.7	1.5	0.2	setosa
4.8	3.4	1.6	0.2	setosa
4.8	3.0	1.4	0.1	setosa
4.3	3.0	1.1	0.1	setosa
5.8	4.0	1.2	0.2	setosa
5.7	4.4	1.5	0.4	setosa
5.4	3.9	1.3	0.4	setosa
5.1	3.5	1.4	0.3	setosa
5.7	3.8	1.7	0.3	setosa
5.1	3.8	1.5	0.3	setosa
5.4	3.4	1.7	0.2	setosa
5.1	3.7	1.5	0.4	setosa

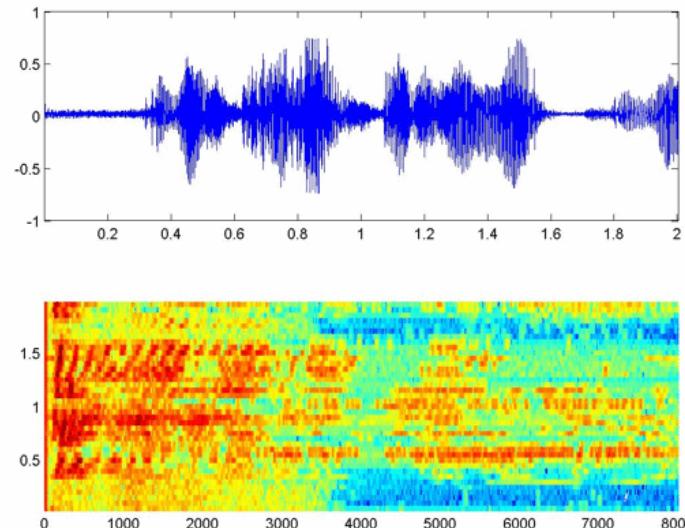
feature vectors

labels

Typical assumption on data:

- i.i.d.
- Surrounding tasks deemed simple(r)

Sequential Supervised Learning



Typical assumptions on data:

- Sequence information matters
- Overall task has many challenging components (e.g., segmentation → recognition → sequence assembly)

Approaches to classifying sequential data

«A bird in the hand...» approach

- Train standard classifier, extend it using a sliding window and post-processing (e.g., smoothing)

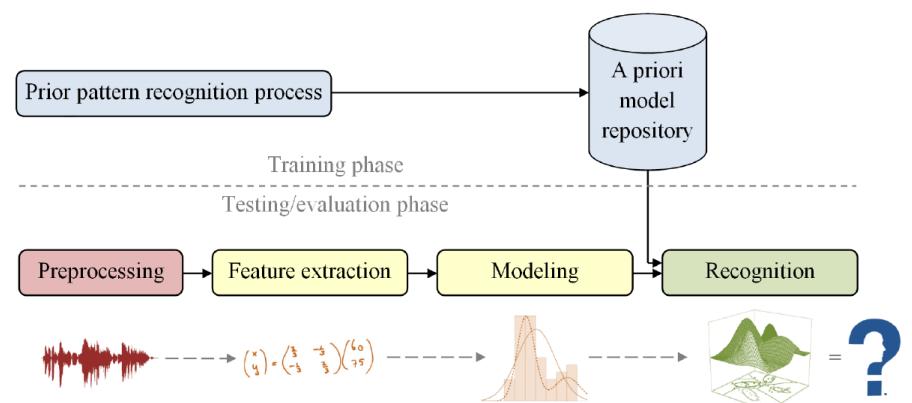
Direct modeling approach

- Train a generative (statistical) model of the sequence generation process (e.g., HMM)

«...two in the bush» approach

- Build a unified pattern recognition processing chain, optimize it globally with a unique criterion

See also: T.G. Dietterich, «Machine Learning for Sequential Data – A Review», 2002



PROPOSED SOLUTION: GLOBAL LEARNING

Images sources: See references on last slide

Example: Reading handwritten strings

Challenge

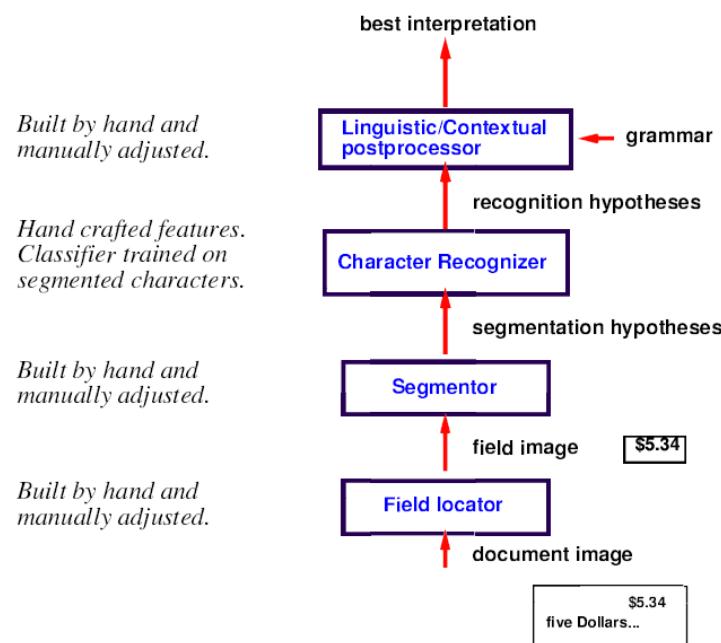
- Identify correct character string («342») on a piece of paper
- Therefore: Find correct segmentation & recognize individual characters



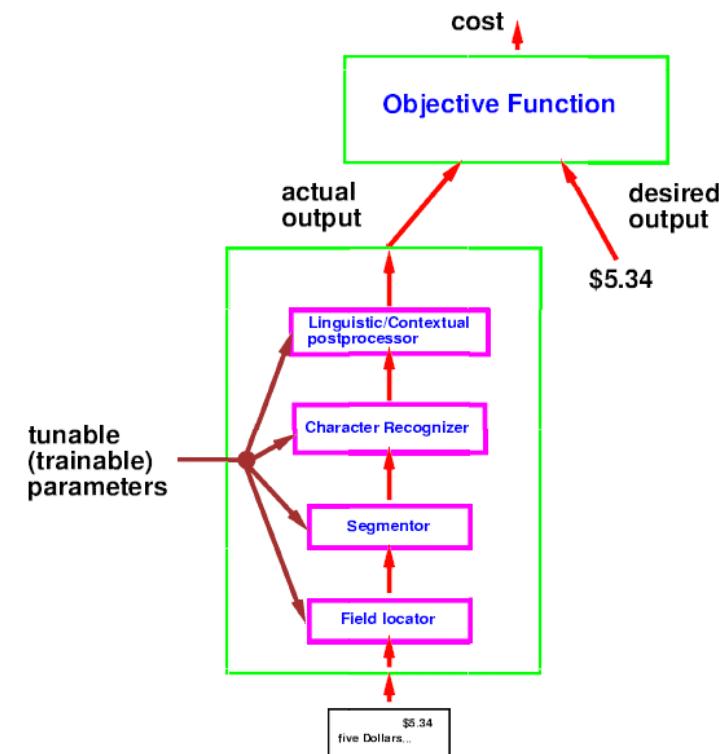
Global Learning

Learning end-to-end

What we know: Traditional pattern recognition system architecture

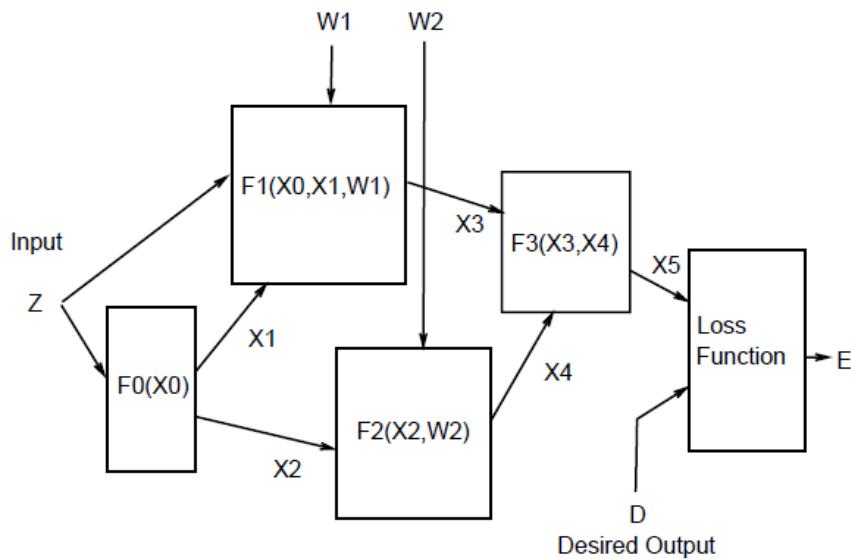


What we want: Train all parameters to optimize a global performance criterion



Foundation: Gradient-based learning

A trainable system composed of heterogeneous modules:



Backpropagation can be used if...

- cost (loss) function is differentiable w.r.t. parameters
- modules are differentiable w.r.t. parameters

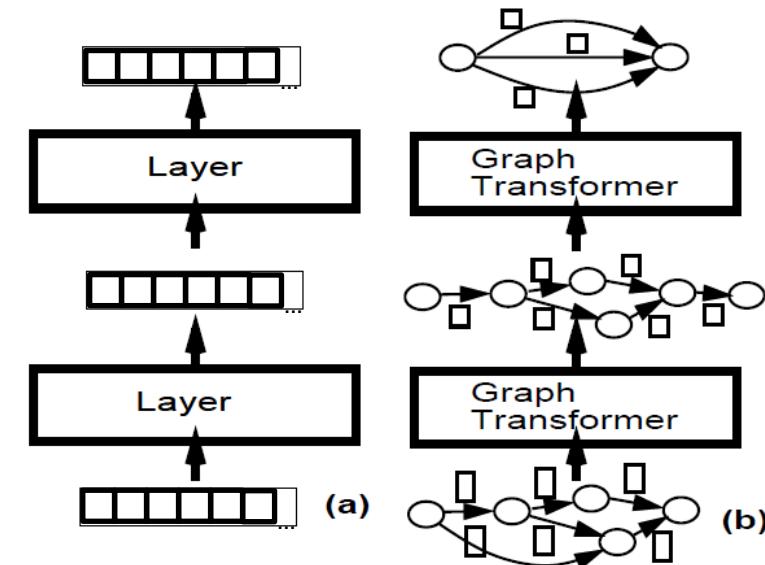
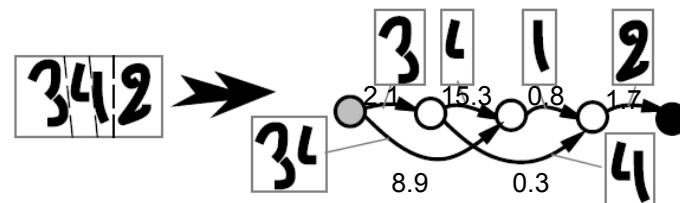
- ➔ Gradient-based learning is the unifying concept behind many machine learning methods
- ➔ Object-oriented design approach:
Each module is a class with a **fprop()** and **bprop()** method

Graph Transformer Networks

Network of pattern recognition modules that successively refine graph representations of the input

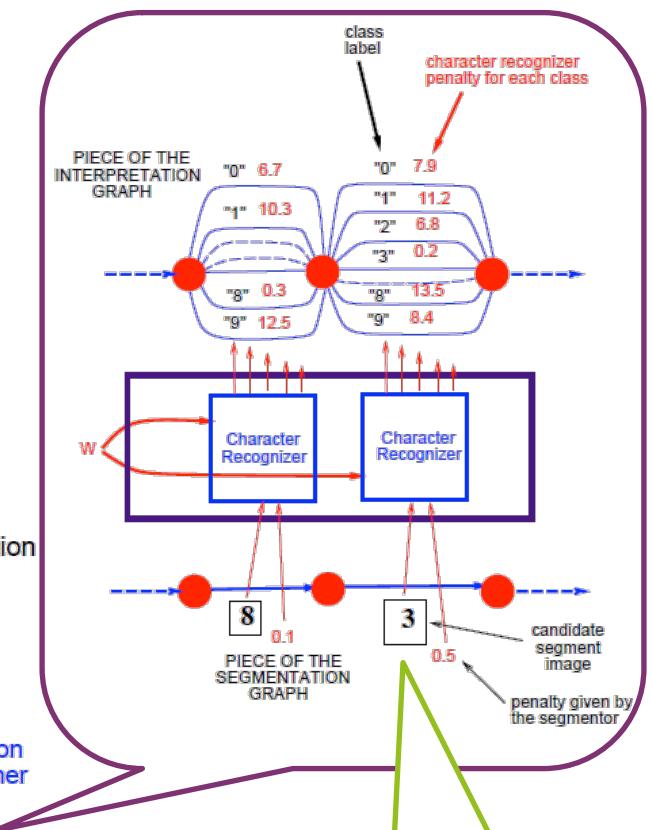
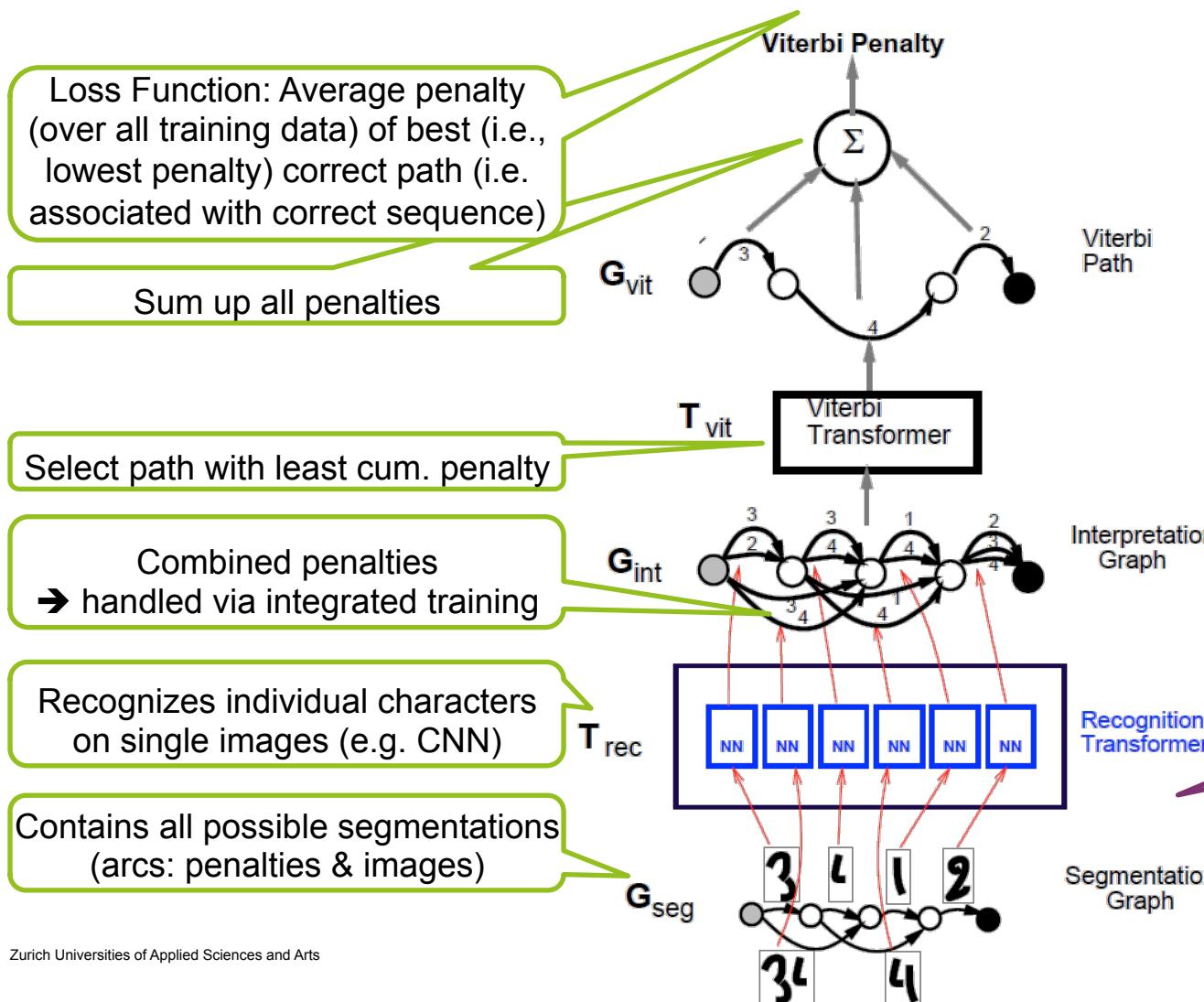
GTNs

- Operate on graphs (b) of the input instead fixed-size vectors (a)
- Graph: DAG with numerical information (“penalties”) at the arcs



→ GTN takes gradients w.r.t. module parameters and numerical data at input arcs

Example: Heuristic oversegmentation ...for reading handwritten strings

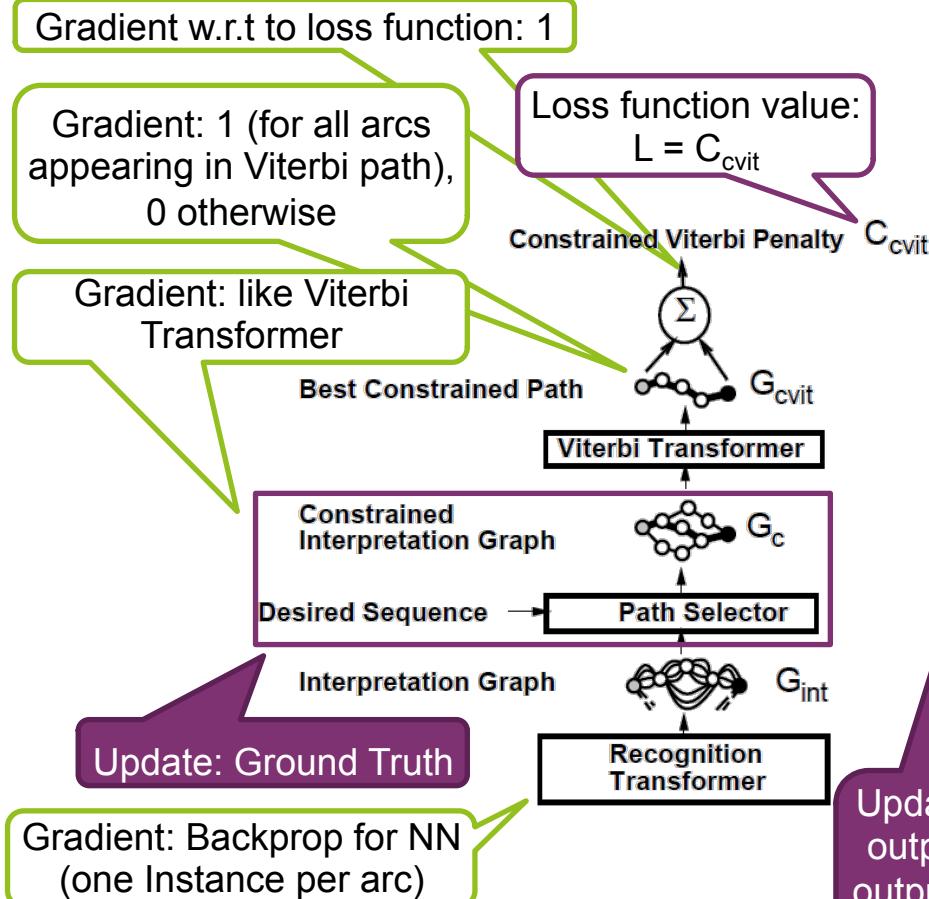


Problems:

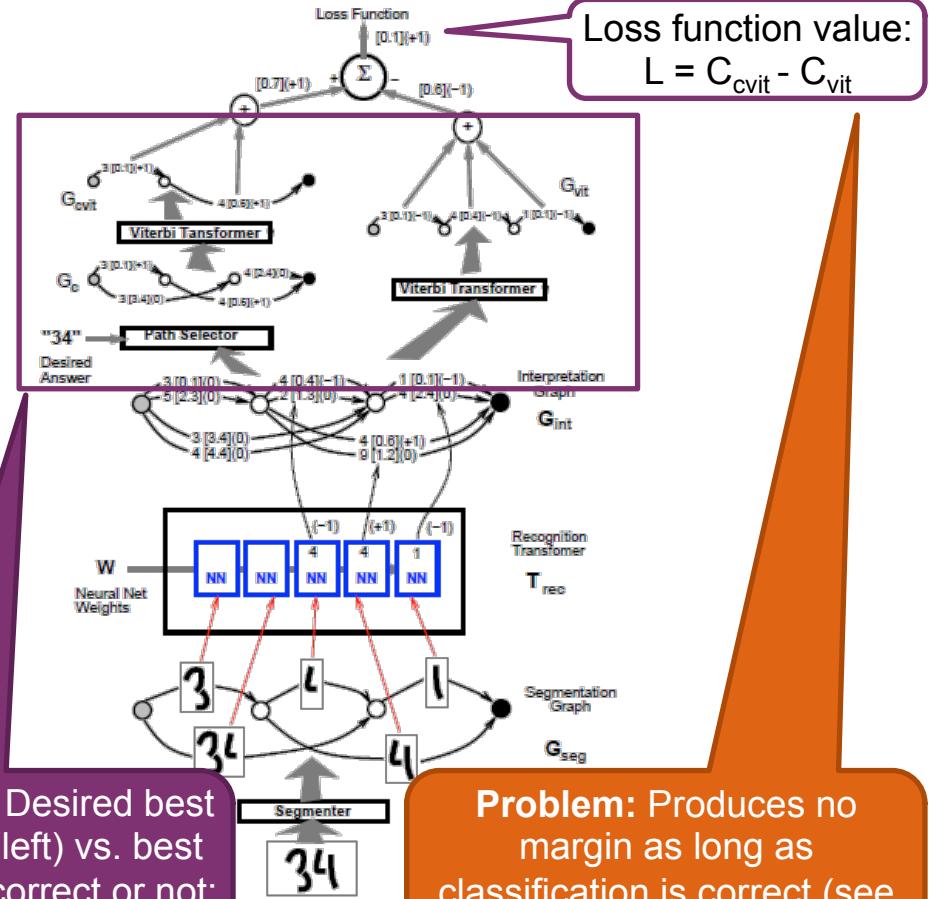
1. Trivial solution possible (Recognizer ignores input & sets all outputs to small values)
2. Penalty does not take competing answers into account (i.e., ignores training signals)

Solved: Discriminative training builds the class-“separating surfaces rather than modeling individual classes independently of each other” → $L=0$ if best path is a correct path.

«Viterbi» training



Discriminative training



Update: Desired best output (left) vs. best output (correct or not; right)

Problem: Produces no margin as long as classification is correct (see paper for solution).

Remarks

Discriminative training

- Uses **all** available training **signals**
- Utilizes “penalties”, **not probabilities**
 - **No need for normalization**
 - Enforcing normalization is “*complex, inefficient, time consuming, ill-conditions the loss function*” [according to paper]
- Is the **easiest/direct way** to achieve the objective of classification (**as opposed to Generative training**, that solves the more complex density estimation task as an intermediary result)

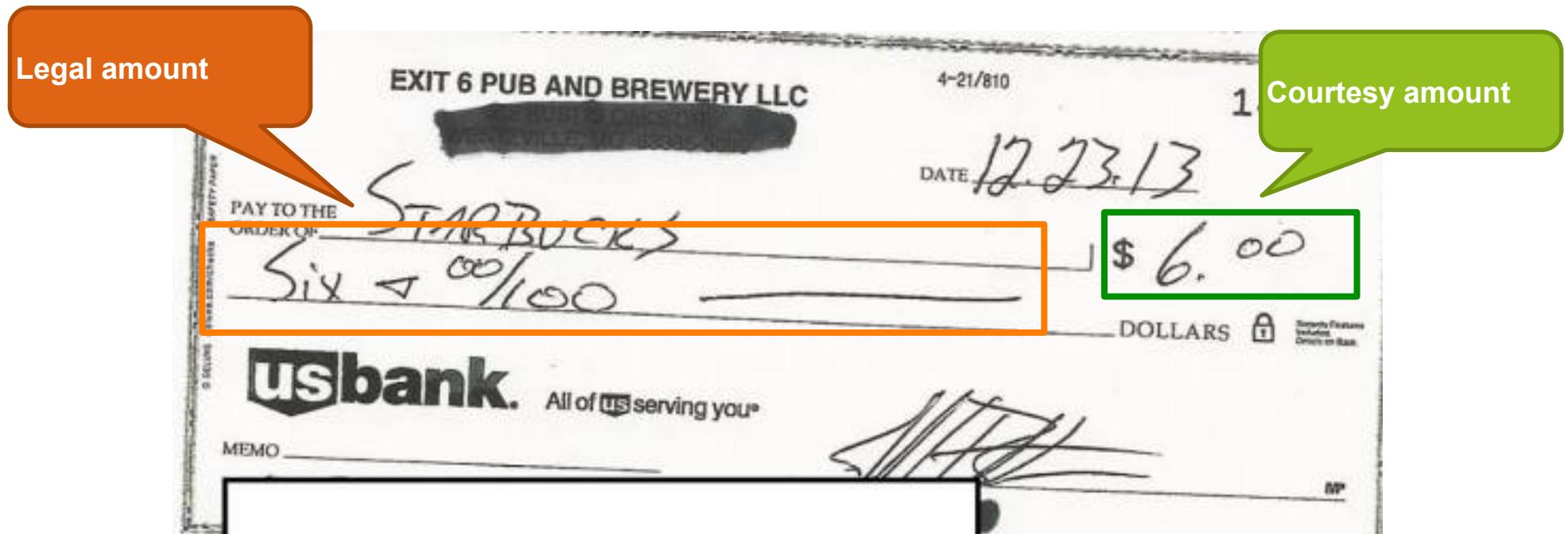
List of possible GT modules

- All **building blocks** of (C)NNs (layers, nonlinearity etc.)
- **Multiplexer** (though not differentiable w.r.t. to switching input)
 - can be used to dynamically rewire GTN architecture per input
- **min**-function (though not differentiable everywhere)
- **Loss** function



EXAMPLE: CHECK READER APPLICATION

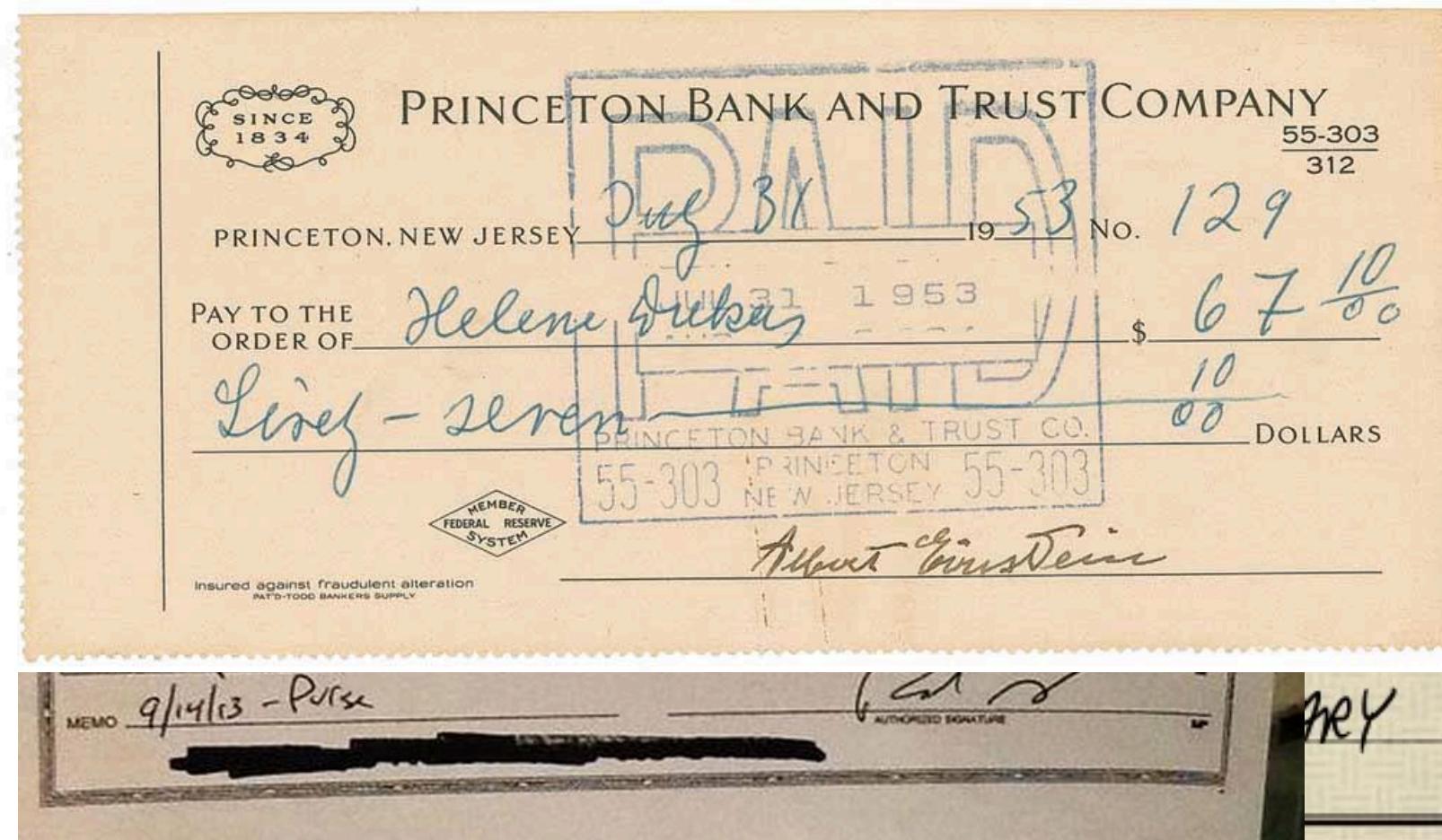
What is the amount?



**easy for human
but time consuming**

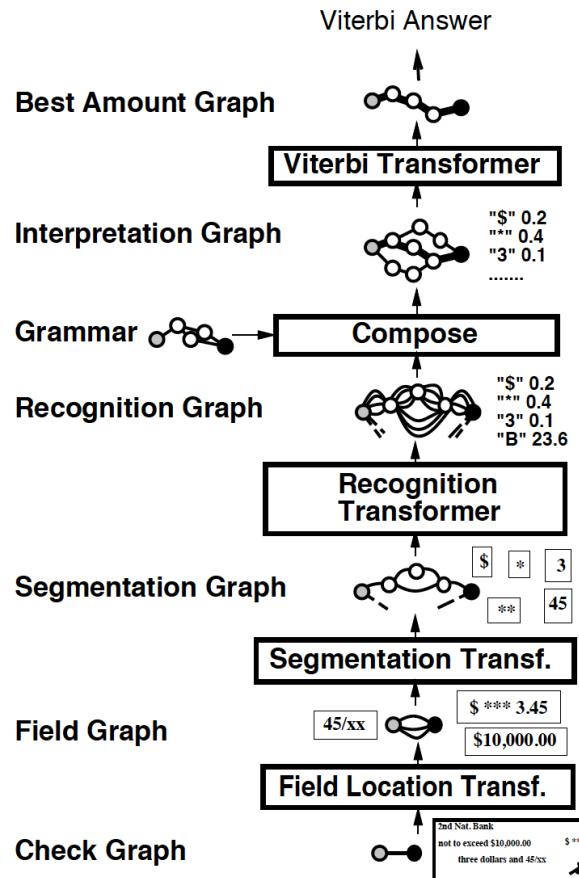
interest in automating the process

A Check Reading System



Goal: Successfully read the amount of \$ 6.00

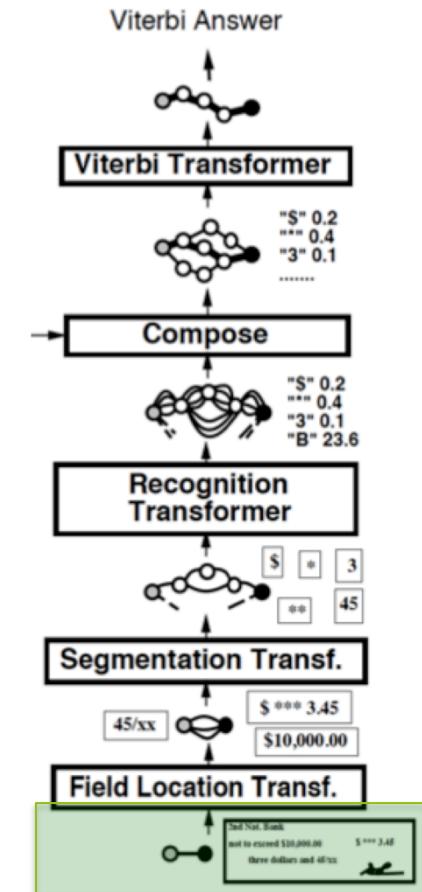
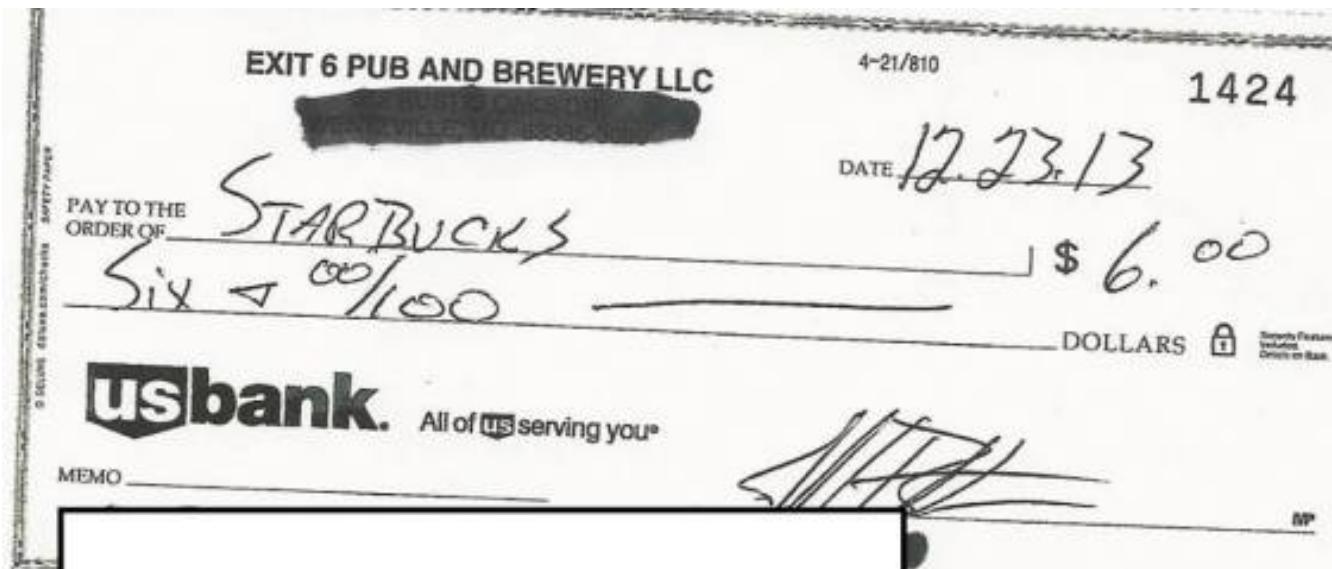
A Check Reading System



- NCR
- Reading millions of checks per day
- GTN-based
- Aimed performance:
 - 50% correct
 - 49% reject
 - 1% error

Check graph

- Input: trivial graph with a single arc that carries the image of the whole check



Field location

Transformer

- Performs classical image analysis (e.g. Connected component analysis, Ink density histogram, Layout analysis)
- Heuristically extract rectangular zones which may contain the amount.

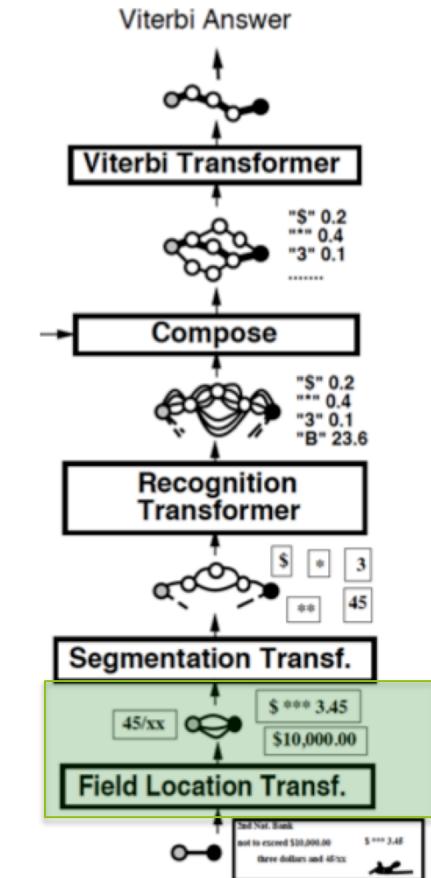
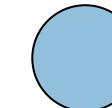
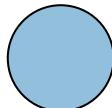
00/100

12.23.13 \$ 6. 00

4-21/810

1424

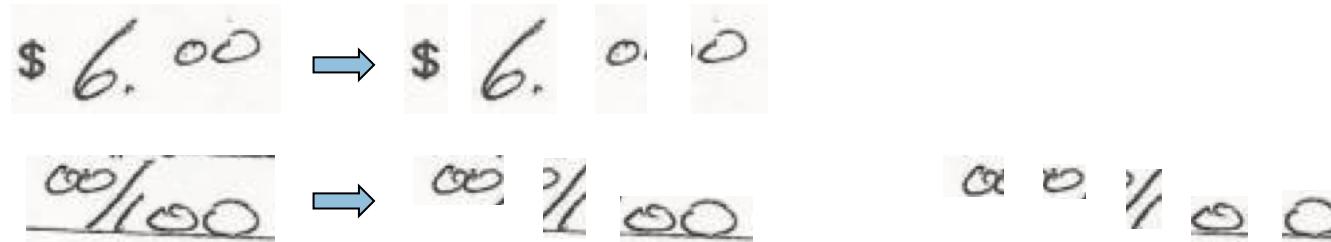
Field Graph



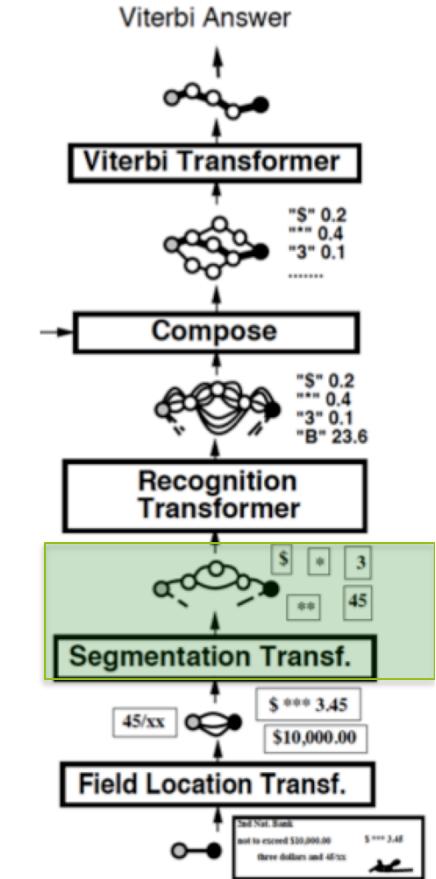
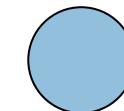
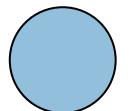
Segmentation

Transformer

- Performs heuristic image processing techniques (here hit an deflect) to find candidate cuts.



Field Graph



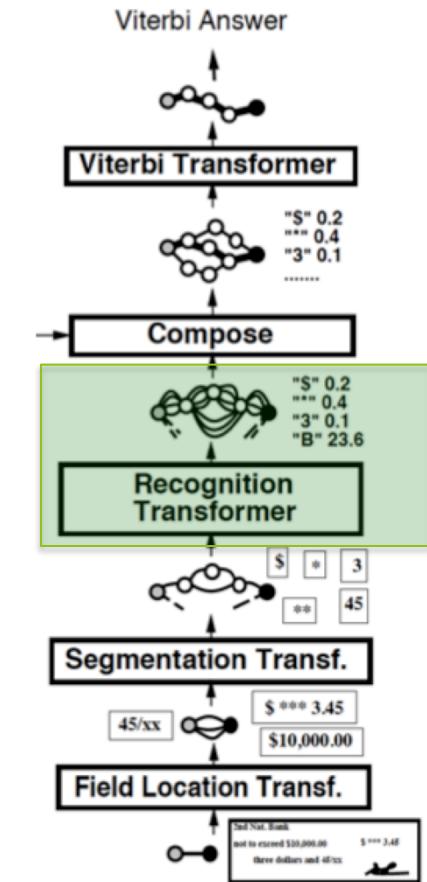
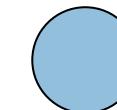
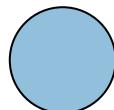
Recognition

Transformer

- Iterates over all arcs in the segmentation graph and runs a character recognizer [here: LeNet5 with 96 classes (ASCII: 95, rubbish: 1)]



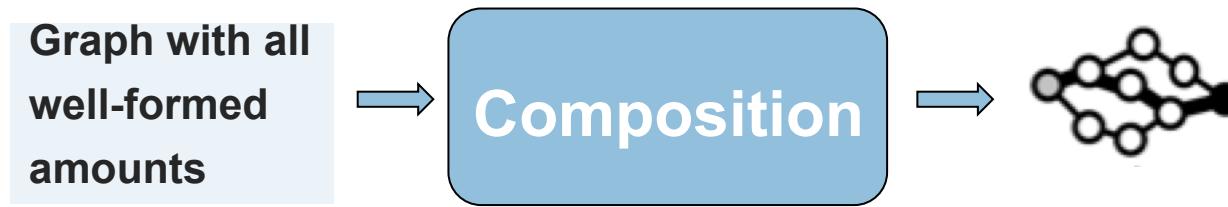
Field Graph



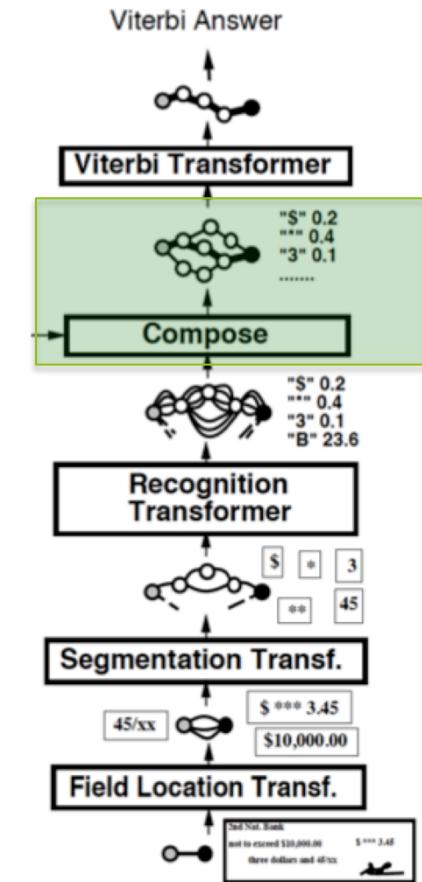
Composition

Transformer

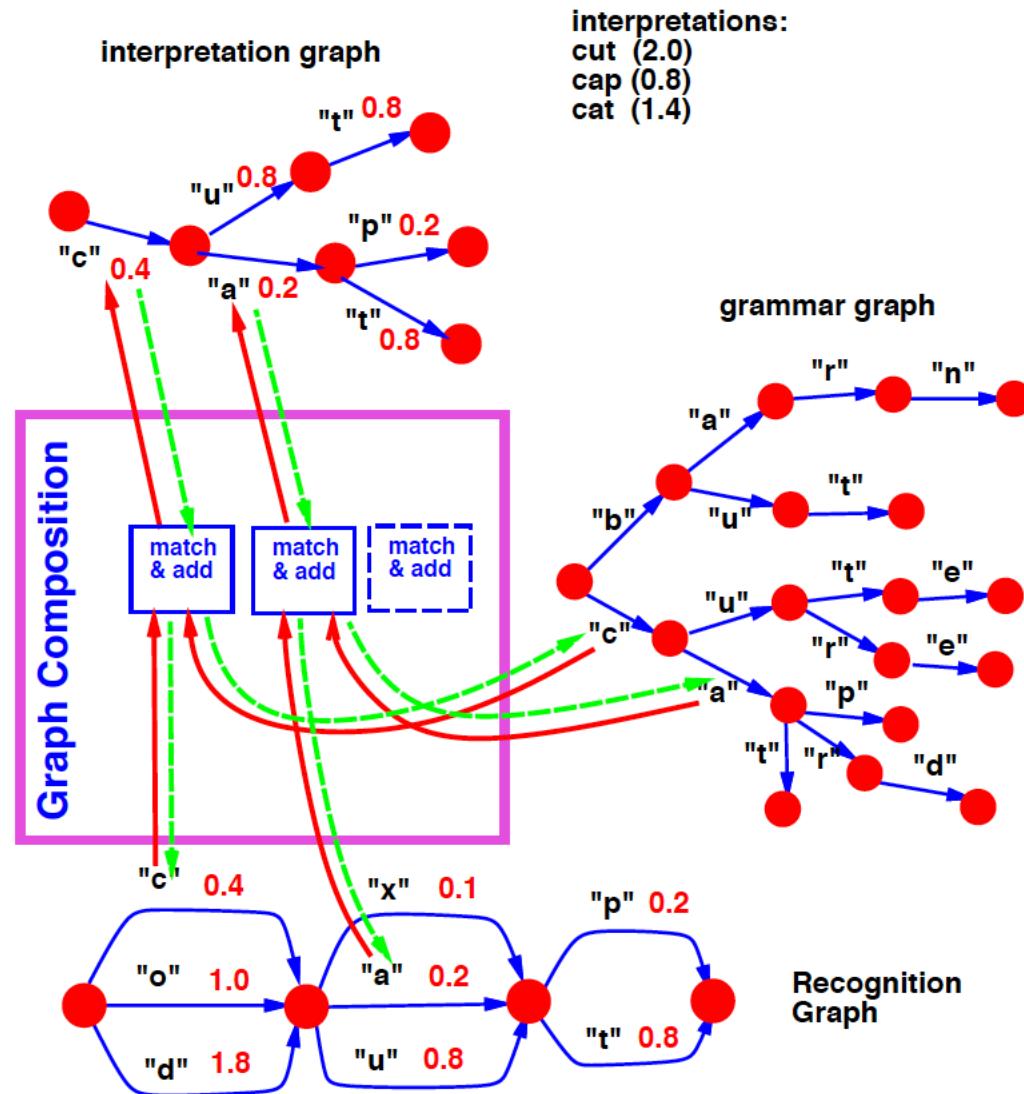
- Takes two graphs, the recognition graph and a grammar graph, to find valid representation of check amounts.



Field Graph



Composition

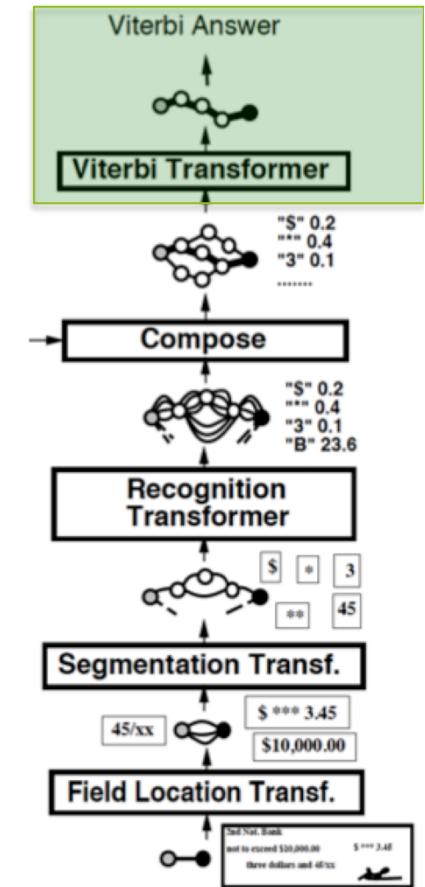
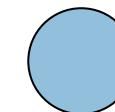
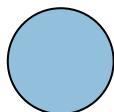


Viterbi

Transformer

- Selects the path with the lowest accumulated penalty.

Field Graph



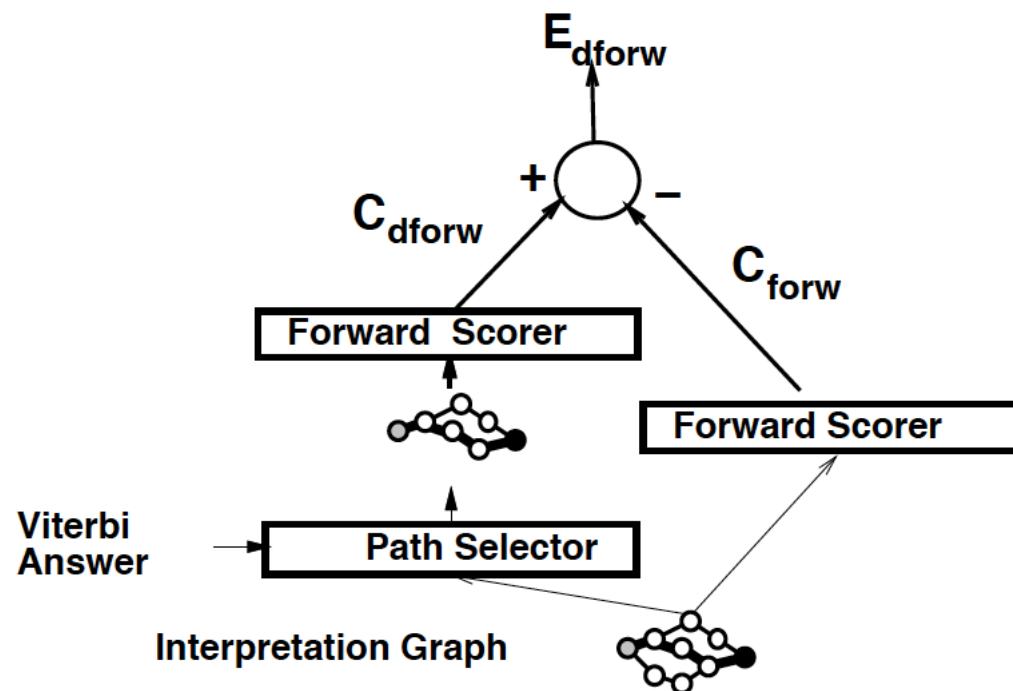
Learning

- Every transformer need to be differentiable
- Each stage of the system has tunable parameters
- Initialized with reasonable values
- Trained with whole checks labeled with the correct amount
- Loss function E minimized by **Discriminative Forward** criterion:

$$L = C_{cvit} - C_{vit}$$

Discriminative Forward

- Difference between:
 - Path with the lowest sum of penalties and the correct Seq
 - Path with the lowest sum of penalties and the allowed Seq



REMARKS / OUTLOOK

Conclusions

Less need for manual labeling

- Ground truth only needed for final result (not for every intermediate result like e.g. segmentation)

Early errors can be adjusted later due to...

- ...unified training of all pattern recognition modules under one regime
- ...postponing hard decisions until the very end

No call upon probability theory for modeling / justification

- Occam's razor: Choose easier discriminative model over generative one
- Vapnik: Don't solve a more complex problem than necessary
- No need for normalization when dealing with "penalties" instead of probabilities → no "other class" examples needed
- Less constraints on system architecture and module selection

Example: Learning to segment without intermediate labels



Zurich University
of Applied Sciences



possible without crop marks?

Gregory To Speak At Coalition Rally

Cavalier Daily Staff Writer
Along with all the events taking place at the University for April 14, Founder's Day, there will be another event that has not been posted on any Sesquicentennial calendar. On the day of Founder's Day, George Jackson, speaker and political agitator at the last national elections, will speak at a rally of the student coalition.

The coalition, contacted Mr. George Jackson, the speaker for the celebration to Mr. Dick Joyce, director of the Cavalier Daily, who said that he could. "I fed soliciting donations from student groups for Mr. Gregory will be contributed to the Transition Program by the Cavalier Daily," said Mr. Joyce.

The student coalition is not the only group that has been invited to Mr. Gregory's speech, but is trying to raise the money to pay for his appearance. "We are going to have a short debate, it was decided by the coalition that it was important to have a 'white and black' debate," said Dick Joyce. "It will be held at Jefferson Hall Auditorium, which is located in Jefferson Hall, the auditorium where the Cavalier Daily is located.

Following the student coalition's rally on April 14, a counter Sesquicentennial celebration is planned. At approximately 2 p.m., the coalition will be gathered in the University Union lobby. After meeting Tuesday night to decide whether or not to give a speech, the coalition members will be gathered in the University Union lobby to attend the "Counter Sesquicentennial".

In addition to the "Counter Sesquicentennial", coalitions members

Further Reading

- Original short paper: Bottou, Bengio & LeCun, “Global Training of Document Processing Systems using Graph Transformer Networks”, 1997
<http://www.iro.umontreal.ca/~lisa/poiteurs/bottou-lecun-bengio-97.pdf>
- Landmark long paper: LeCun, Bottou, Bengio & Haffner, “Gradient-Based Learning Applied to Document Recognition”, 1998
<http://yann.lecun.com/exdb/publis/pdf/lecun-01a.pdf>
- Slide set by the original authors: Bottou, “Graph Transformer Networks”, 2001
<http://leon.bottou.org/talks/gtn>
- Overview: Dietterich, “Machine Learning for Sequential Data: A Review”, 2002
<http://eecs.oregonstate.edu/~tgd/publications/mlsd-ssspr.pdf>
- Recent work: Collobert, “Deep Learning for Efficient Discriminative Parsing”, 2011
http://ronan.collobert.com/pub/matos/2011_parsing_aistats.pdf

