



Data Compression With and Without Deep Probabilistic Models

Prof. Robert Bamler

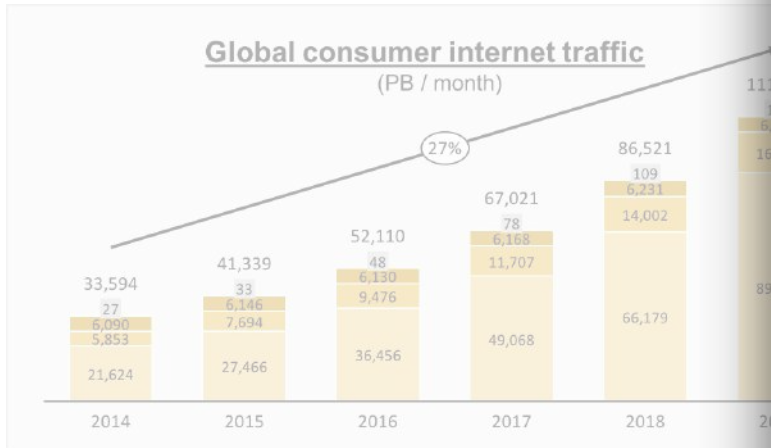
Lecture #1 of course “Data Compression With Deep Probabilistic Models”

University of Tuebingen • 21 April 2022



Why Are We Here?

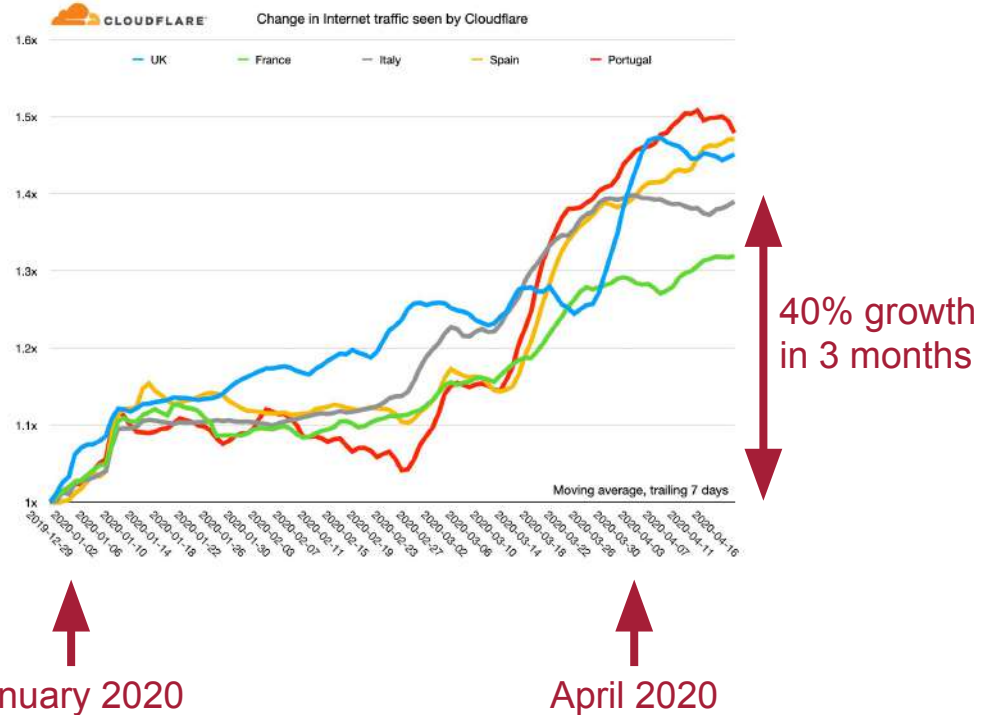
Projection by Cisco from 2015:
27% growth per year



Source: Cisco VNI, 2015

[<https://www.aploris.com/blog/charts/stack-bar-chart-examining-the-growth-of-global-consumer-internet-traffic/>]

Observation by Cloudflare in 2020



[<https://blog.cloudflare.com/recent-trends-in-internet-traffic/>]



a) common data types

NETFLIX

 **YouTube**



zoom

The screenshot shows the InfoQ website interface. At the top, there's a navigation bar with 'InfoQ' logo, a user profile icon, and a notification bell with a '1' badge. Below the navigation bar are several category tabs: 'Streaming', 'Machine Learning', 'Reactive', 'Microservices', and 'Container'. A blue pill-shaped tag reads 'AI, ML & DATA ENGINEERING'. A yellow highlight box contains the text 'QCon Plus (May 17-28): Stay Ahead of Emerging Software Trends'. The main article title is 'NVIDIA's AI Reduces Video Streaming Bandwidth Consumption by 10x'. Below the title are social interaction buttons for 'LIKE', 'DISCUSS', and a bookmark icon. The article date is 'OCT 20, 2020' and the estimated reading time is '3 MIN READ'. The author is 'Anthony Alford' with a 'FOLLOW' button. The article text begins with 'GPU-manufacturer NVIDIA announced their Maxine platform for AI-enhanced video-conferencing services, which includes a technology that can reduce bandwidth requirements by an order of magnitude. By moving much of the data' and ends with a URL: '[https://www.infoq.com/news/2020/10/nvidia-video-bandwidth/]



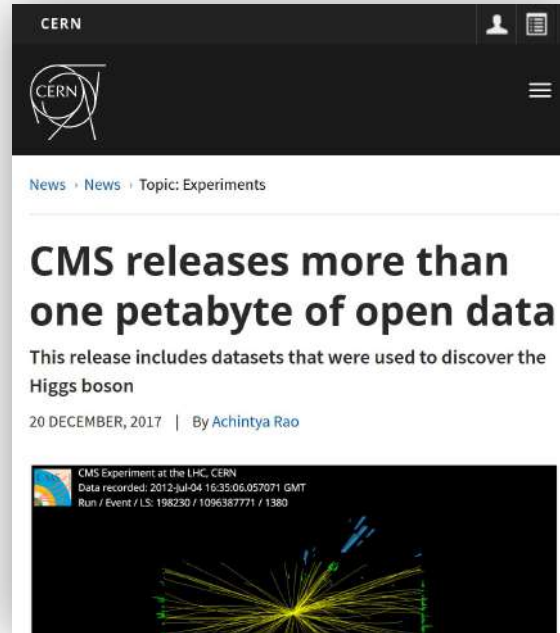
a) common data types

NETFLIX



zoom

b) specialized data types



[home.cern]



[openneuro.org]



From information theory to applications:

- theory of communication (*aka* information theory)
- theoretical bounds for lossless & lossy compression
- generic compression algorithms (“entropy coders”)
 - mathematical proofs of optimality & practical implementations
- probabilistic models & probabilistic inference
 - mathematical derivations & practical implementations
- interplay between (deep) probabilistic models & compression algorithms

Administrative Issues





I Need Your Feedback

This is an advanced course. Please don't be shy and ask questions.



<https://robamler.github.io/teaching/compress22/>

Contains:

- dates & times
- slides & notes (including these slides)
- problem sets (including code) & solutions
- additional video material (occasionally)
- links to moodle (which has zoom link)



Lecture Times & Tutorials

Lectures: Thursdays, 12:15 - 13:45, room TTR2 (Maria-von-Linden-Straße 6)

- mix of whiteboard, slides, and some brief programming demos

Tutorials: Fridays, 12:15 - 13:45, room A302 (Sand 1)

- Problem sets online after each Thursday lecture.
- Tutorials will discuss the problem set from 8 days ago (except tomorrow, where we'll discuss today's introductory problem set).
- Problem sets are not graded but there will be anonymous self-evaluation polls on moodle. **Please participate in your own interest!**



Grading

- 6 ECTS
- exam date TBD (any conflicts?)

Getting to Know You

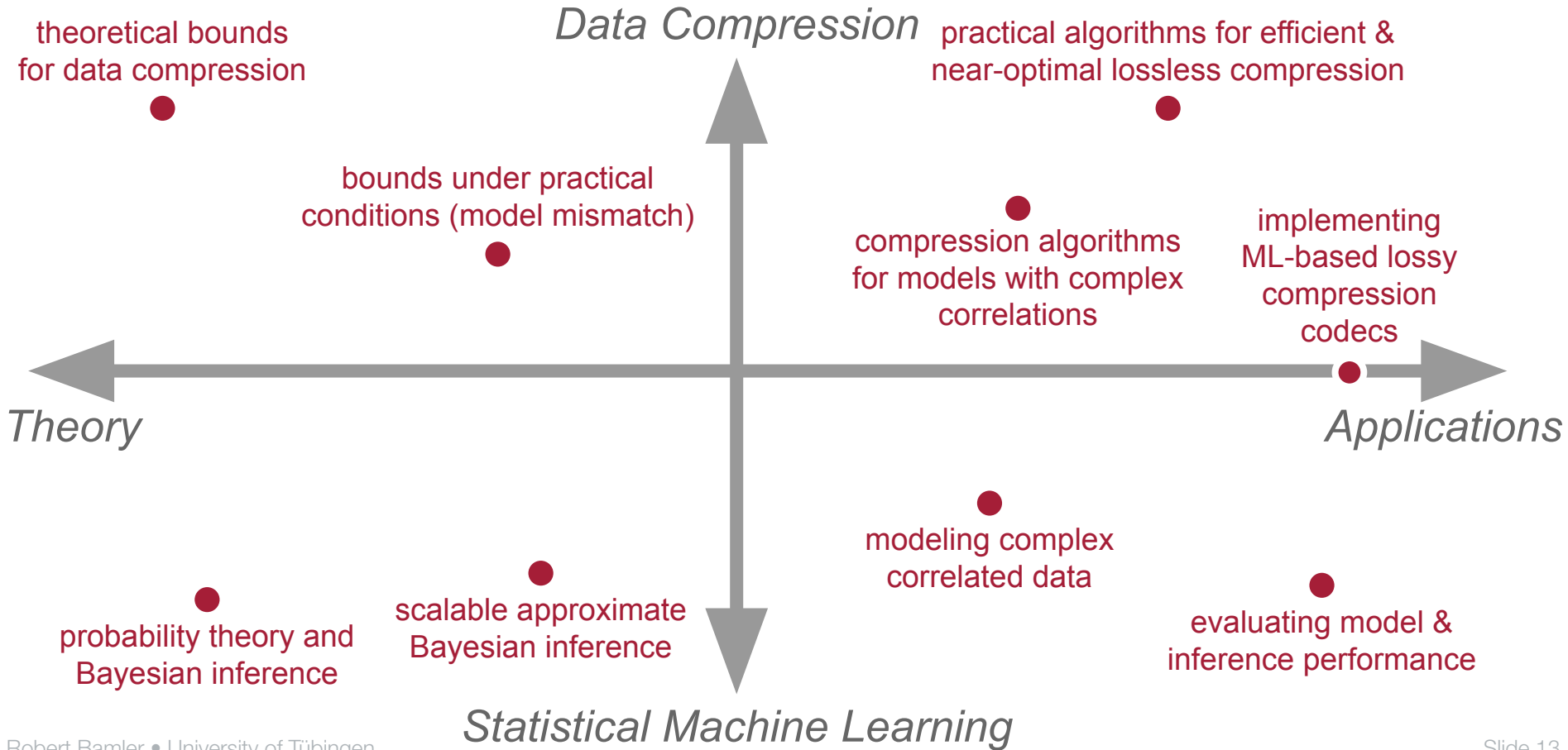


Course Overview



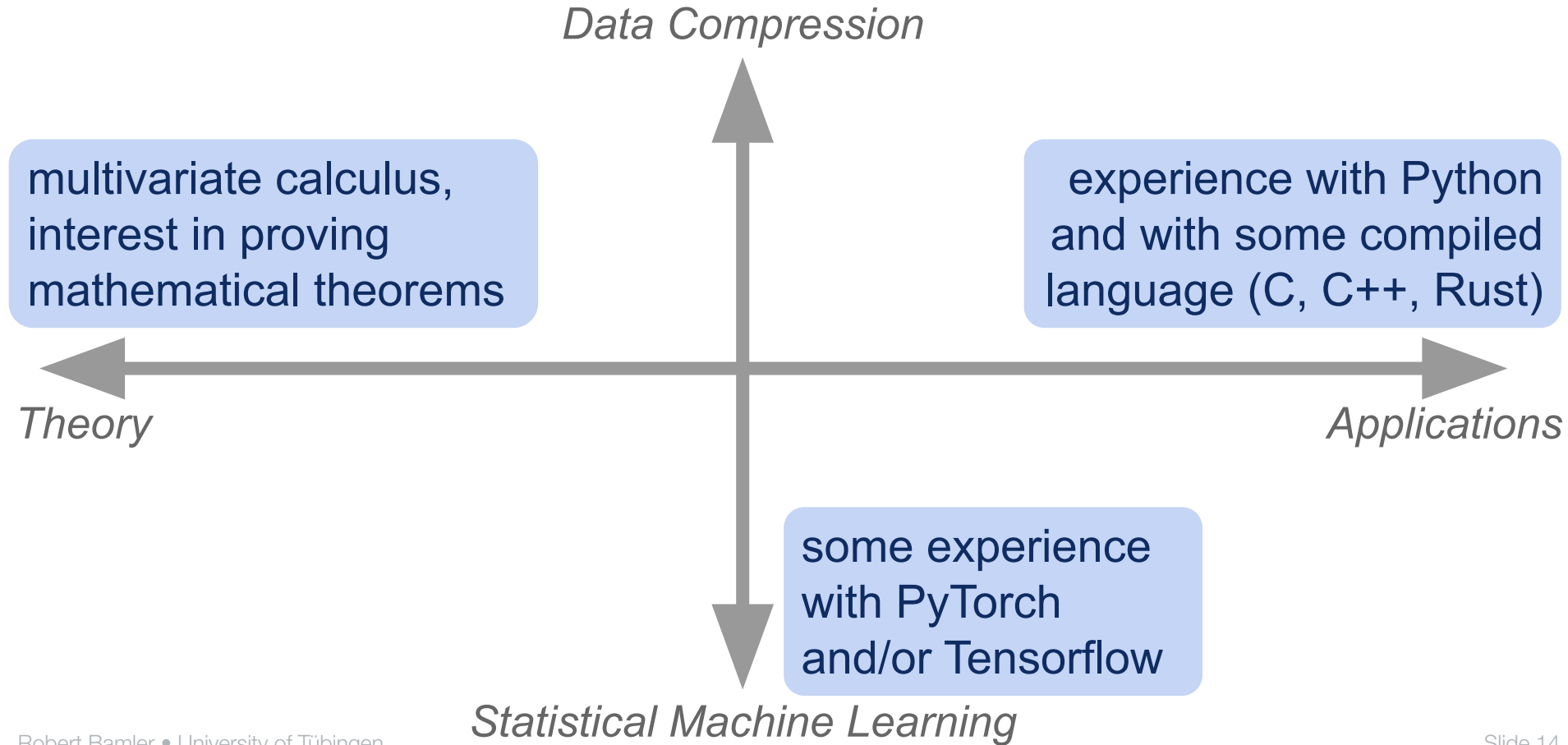


Spectrum of Topics





Prerequisites





Agenda

<https://robamler.github.io/teaching/compress22/>



Books:

- **Information Theory:** MacKay, *Information theory, inference and learning algorithms*. Cambridge university press, 2003
legal free PDF: <http://www.inference.org.uk/mackay/itprnn/book.html>
- **Probabilistic ML:** Murphy, Kevin P. *Machine learning: a probabilistic perspective*. MIT press, 2012.

Videos:

- lectures by David MacKay:
<https://youtube.com/playlist?list=PLruBu5BI5n4aFpG32iMbdWoRVAA-Vcso6>
- mathematicalmonk: <https://youtube.com/playlist?list=PLE125425EC837021F>
- videos from last year's course: see link on course website

Related Lecture:

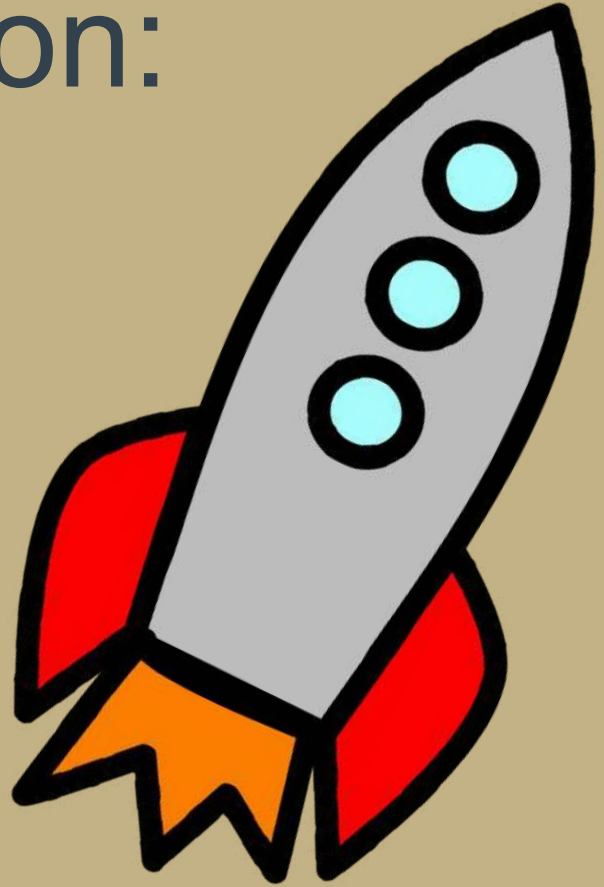
- Probabilistic Machine Learning by Prof. Hennig (this term)



USP: How is This Lecture Different?

- full spectrum from theory to applications
 - we'll *prove* theorems but we'll also *implement* real compression codecs
- full spectrum from information theory to statistical machine learning
 - “models and algorithms”
- highly active field of research
 - lots of recent developments not yet covered in secondary literature
 - both compression algorithms & deep probabilistic models
- interactivity
 - mini-problems in class so I can assess which topics I need to clarify

Neural Compression:
a highly active field
of research.





Active Research Topics (Selection)

- entropy models for important data types (especially video)
 - normalizing flows
 - implicit generative models
- capturing correlations in compression algorithms
 - compression algorithms for hierarchical models
 - compression algorithms for stochastic sampling
- quantization of real-valued neural activations
- model compression
- semantic compression
- learned distortion measures
- computational efficiency of neural compression methods

collaborators:



Yibo Yang
(UC Irvine)

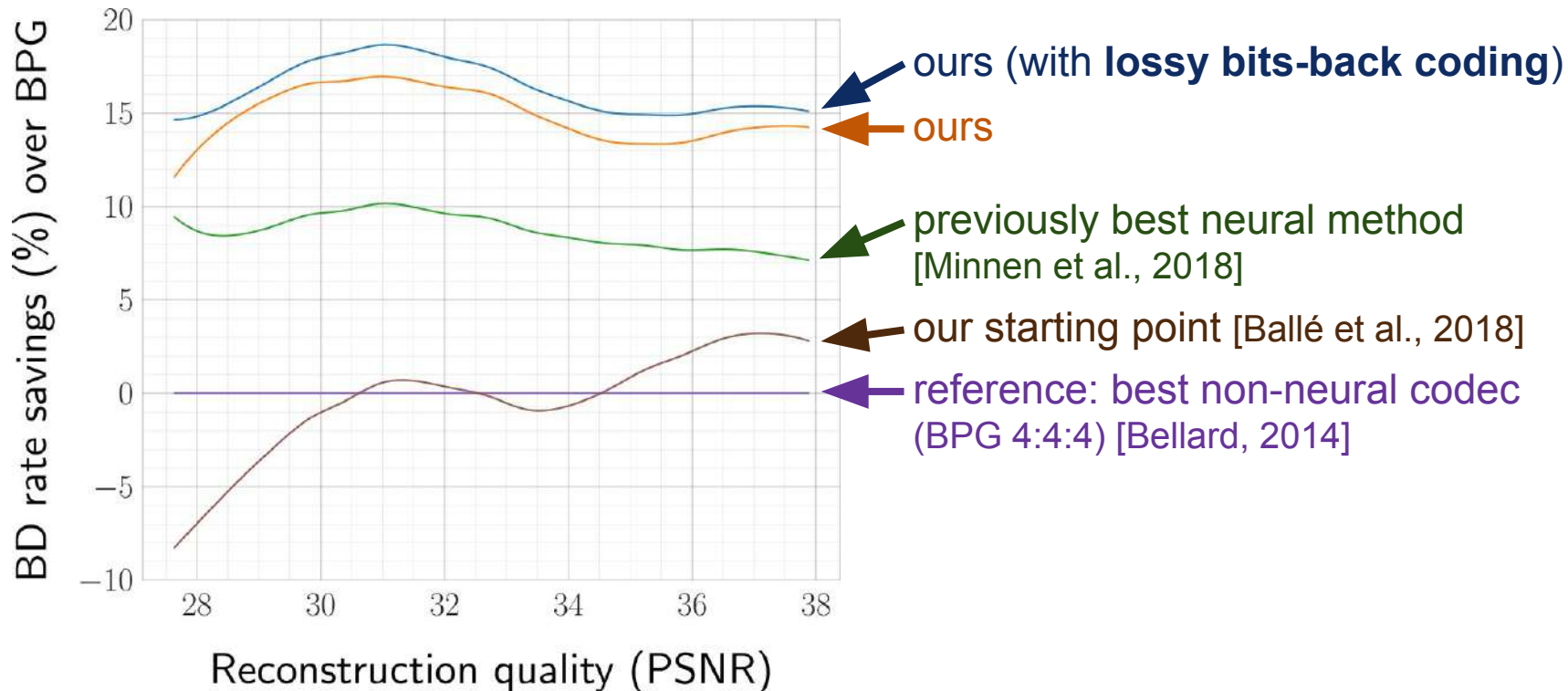


Stephan Mandt
(UC Irvine)



Example: Lossy Image Compression

[Yang, RB, Mandt, NeurIPS 2020]



What's the Population of Rome?

▶ On 30 April 2018: 2,879,728

▶ In the year 500 AD: 100,000



original



JPEG @ 0.24 bits per pixel



Ours @ 0.24 bits per pixel





Case Study: Student Project

Neural compression method
running *in a web browser*

© Alexander Conzelmann
and Shrisha Bharadwaj

ZIPNET


Zipnet - a neural image codec entirely in your browser!

Try it out with your own image. A trained neural network will find an efficient encoding for your image and you can then download it. The encoding will be lossy (similar to JPEG - but it often looks more natural to the eye). You can then re-upload the encoded image to see the reconstruction. And the best thing about it - this algorithm runs entirely in your browser. No servers or anything involved! You can try it out below. If you need additional instructions or want explanations on what you see refer to the [Detailed usage](#).

ENCODING IMAGE

Choose a **.png** image you want to encode (max size 256x256).

Icon-128.png

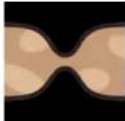
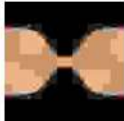


Previous size: 4000 byte, Encoded size: 572 byte (14.30 % of original), 0.28 bpp, Encoding time: 232 ms

DECODING IMAGE

Upload a Zipnet encoded image. The file should be called something like "zipnet-enc-timestamp.bin".

zipnet-enc-1850533815702.bin

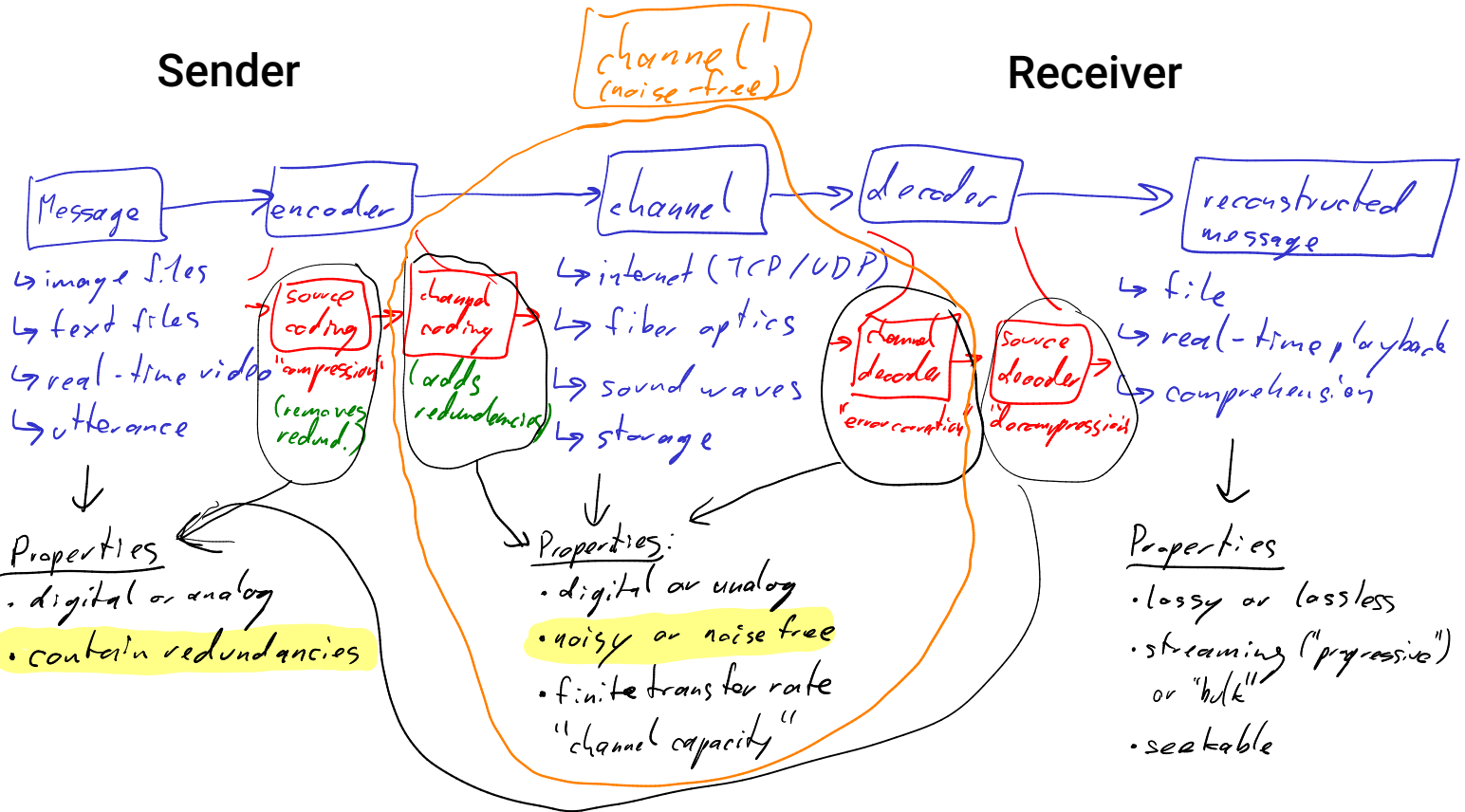
Zipnet:  JPEG: 

Decoding time: 240ms
Zipnet: 0.28 bpp 36.58 PSNR
JPEG: 0.55 bpp 23.76 PSNR

Data Compression With and Without Deep Probabilistic Models

Lecture 1 (21 April 2022)

Problem Setting: communication over a channel



"Source-channel separation theorem" → we'll cover this in more detail when we talk about lossy compression

Goal: transmit message from Sender to Receiver:

- fast, i.e., using the channel as little as possible
- reliably, i.e., without errors or with as little (relevant) distortion as possible

Question 1: Assume you want to transmit some message over some channel. The message is given as a string of N bits, and the channel transmits one bit at a time but it occasionally introduces a random bit flip. How many bits do you have to transmit over the channel if you want to be certain beyond reasonable doubt that the receiver of the message can decode it without errors?

- A) N bits
- B) more than N bits
- C) fewer than N bits

D) It depends

- If we have any prior knowledge then we can compress the message (e.g., if we know that the message is an ASCII representation of English text then we know that the letter "e" is most frequent and that the letter "q" is almost always followed by a "u").
- Error correction requires to add some bits
- see Questions 2 & 3 below for further thoughts

Question 2: Assume you have a noise-free channel and a message that contains some redundancies (e.g., English text). What should the encoder and decoder do with these redundancies if our goal is efficient communication (i.e., using the channel as little as possible)?

- encoder should remove redundancies
 - decoder should reconstruct the redundancies
- } compression

E.g., construct a code book, on which encoder & decoder have to agree beforehand

Question 3: Now let's assume that the channel introduces noise (e.g., occasional bit flips). What additional task do encoder and decoder have to do? Think again about redundancies.

- introduce new redundancies that are designed for the channel

Lossless Compression I: Symbol Codes

Problem Setting

- goal: efficient lossless communication over a noise free channel
- sender has some message \underline{x} wants to transmit it losslessly to a receiver in as few bits as possible

We assume that the message \underline{x} is a sequence of symbols from a discrete alphabet: \mathcal{X} is finite or countably infinite

$$\underline{x} = (x_1, x_2, x_3, \dots, x_k) \equiv (x_i)_{i=1}^k \quad \text{where } x_i \in \mathcal{X} \quad \forall i$$

\uparrow
alphabet

Thus, our goal is to find a mapping:

$$\underline{x} \mapsto \text{bit string} \in \{0, 1\}^* \quad \leftarrow \text{Kleene star}$$

such that:

- injective (i.e. invertible)
- bit strings are short

Symbol Codes: Map each symbol in the message to a bit string, then simply concatenate these bit strings

\leftarrow or $\{0, 1, \dots, B\}$ (with $B \geq 2$)

• code book $C: \mathcal{X} \rightarrow \{0, 1\}^*$

• to encode a message: concatenate

$$C^*(\underline{x}) = C^*(x_1, x_2, \dots, x_k) := C(x_1) \parallel C(x_2) \parallel \dots \parallel C(x_k)$$

\downarrow
concatenation

Some nomenclature:

- C : "Code Book"
- C^* : "code"
- $C(x)$: "code word of symbol x "
- Def: $l(x) :=$ the length of $C(x)$ (in bits)

Examples of Symbol Codes

1) Morse Code $B = 3$ (dot, dash, pause)

International Morse Code

1. The length of a dot is one unit.
2. A dash is three units.
3. The space between parts of the same letter is one unit.
4. The space between letters is three units.
5. The space between words is seven units.

A	• —	U	• • —
B	— • • •	V	• • • —
C	— • — •	W	• — • —
D	— • •	X	— • • —
E	•	Y	— • — —
F	• • — •	Z	— — • •
G	— — •		
H	• • • •	1	• — — — —
I	• •	2	• • — — —
J	• — — — —	3	• • • — —
K	— • — — —	4	• • • • —
L	• • • •	5	• • • • •
M	— —	6	— • • • •
N	— •	7	— — • • •
O	— — —	8	— — — • •
P	• — — • •	9	— — — — •
Q	— — • — —	0	— — — — —
R	• — • •		
S	• • •		
T	—		

2) UTF-8 $B = 256$

3) "Simplified Game of Monopoly":

- throw a pair of dice and record the sum of their results as a symbol
- repeat this process several times
- for simplicity, let's use (fair) 3-sided dice:

$$\mathcal{X} = \{2, 3, 4, 5, 6\}$$

x	$p(x)$	$C^{(1)}(x)$	$C^{(2)}(x)$	$C^{(5)}(x)$
2	$1/9$	10	010	010
3	$2/9$	11	011	01
4	$1/3$	100	100	00
5	$2/9$	101	101	11
6	$1/9$	110	110	110

↑
problem; $C^{(1)*}$ is not injective

↑
claim: $C^{(5)*}$ is injective
• an optimal code for \mathcal{X}