# DATA COMPRESSION WITH DEEP PROBABILISTIC MODELS
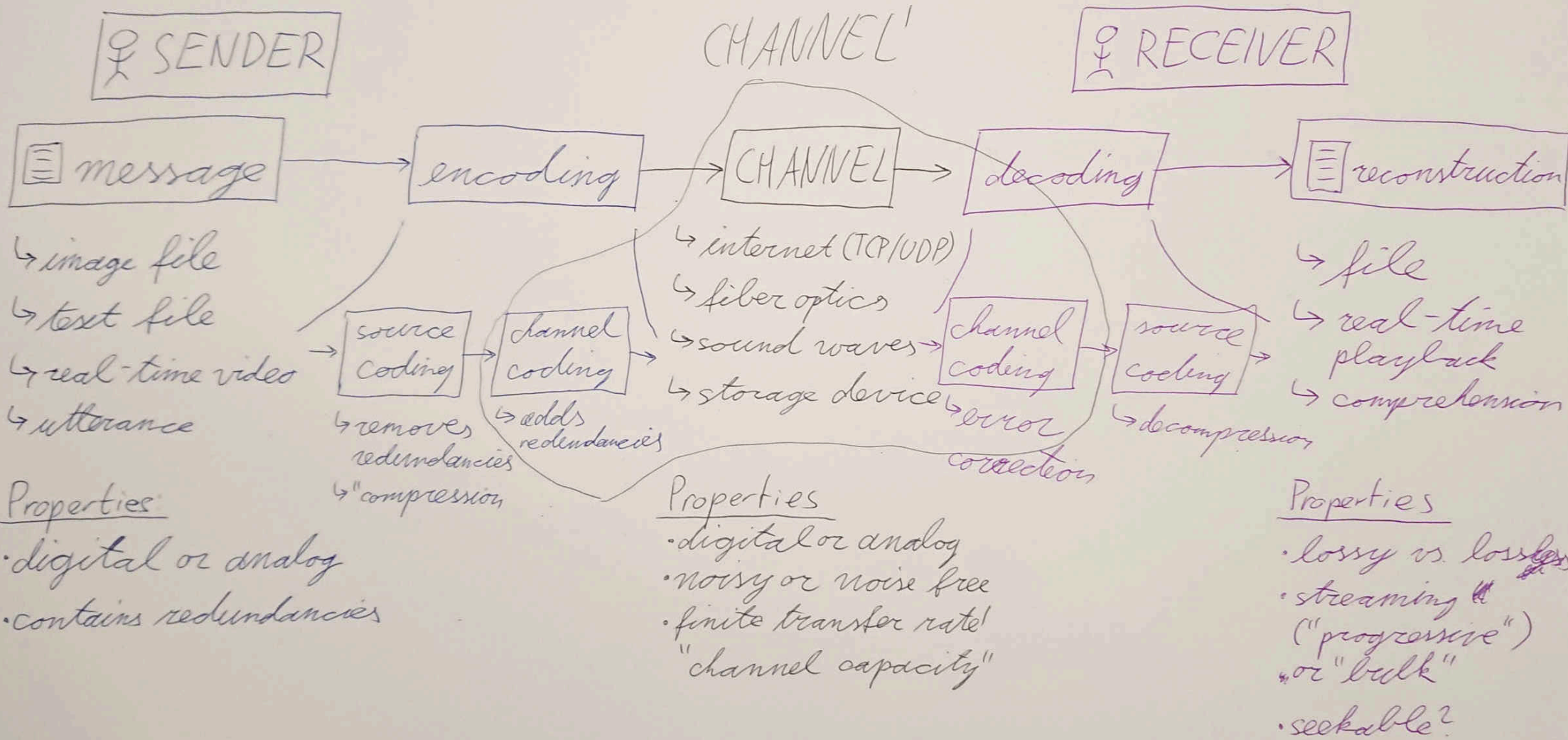
Robert Bamler, University of Tuebingen

## Problem Setting: Communication over a channel

**⚲ SENDER**     CHANNEL'     **⚲ RECEIVER**

message → encoding → CHANNEL → decoding → reconstruction

message:
- image file
- text file
- real-time video
- utterance

encoding:
- source coding
  - removes redundancies
  - "compression
- channel coding
  - adds redundancies

CHANNEL:
- internet (TCP/UDP)
- fiber optics
- sound waves
- storage device

decoding:
- channel coding
  - error correction
- source coding
  - decompression

reconstruction:
- file
- real-time playback
- comprehension

Properties:
- digital or analog
- contains redundancies

Properties:
- digital or analog
- noisy or noise free
- finite transfer rate "channel capacity"

Properties:
- lossy vs. lossless
- streaming ("progressive")
- „or "bulk"
- seekable?

Goal: transmit message from S to R: fast + reliable

# LOSSLESS COMRESSION I: SYMBOL CODES

## Problem Setting

- communicate over a noise free channel
- sender has message $\underline{x}$, wants to transmit it losslessly to receiver in as few bits as possible
- encoder:

$$\underline{x} \longmapsto C^*(\underline{x}) \in \{0, 1\}^* \quad \leftarrow \text{Kleene star}$$

$\underbrace{\{0,1\}^*}$ set of all bit strings of arbitrary length

- more generally: $C^*(\underline{x}) \in \{0, ..., B-1\}^*$ with $B \in \{2, 3, 4, ...\}$ ("B-ary code") (commonly: $B=2$)

## Symbol Codes

- message $\underline{x}$ is a sequence of symbols $x_i$ from a discrete alphabet $\mathcal{X}$:

$$\underline{x} = (x_1, x_2, ..., x_k) \equiv (x_i)_{i=1}^k$$

where $x_i \in \mathcal{X} \; \forall i$ and $k \in \mathbb{N}$ and $\mathcal{X}$ is finite (or countably infinite)

- encoder: $C^*(\underline{x}) = C(x_1) \| C(x_2) \| ... \| C(x_k)$

$\underset{\uparrow \quad \uparrow \quad \uparrow}{\text{concatenation}}$

↳ $C$ is called the "code book"
↳ $C(x)$ is called the "code word for symbol $x \in \mathcal{X}$"
↳ Def: $\ell(x) := $ length of $C(x)$ (i.e., number of bits)

# Examples of Symbol Codes

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

2) UTF-8: ($B=2$)
   - $\mathcal{X} = \{$all UNICODE code points$\}$
   - $C(x) = $ UTF-8 representation of $x$
   - $\ell(x) \in \{8, 16, 24, 32\}$   (bits)

3) "Simplified game of Monopoly":
   - throw a pair of dice several times, after each time, write down their sum as a new symbol $x$;
   - for simplicity, let's use 3-sided dice

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

   $\underset{1+1}{\uparrow} \quad \underset{1+2,}{\uparrow} \qquad\qquad \underset{3+3}{\uparrow}$

   $\underset{2+1}{} \quad \boxed{C^*((2,6)) = \underline{10}\,\underline{110} = C^*((5,2))}$

· possible code books:
   - $C^{(1)}(x) = $ binary representation of $x$
   - $C^{(2)}(x) = \underline{\quad\quad} " \underline{\quad\quad} (x-2)$
   - $C^{(3)}(x) = \underline{\quad\quad} " \underline{\quad\quad} (x-2)$,
     padded to consistent length'
   - $C^{(4)}(x), C^{(5)}(x)$: see table

| $x$ | $C^{(1)}(x)$ | $C^{(2)}(x)$ | $C^{(3)}(x)$ | $C^{(4)}(x)$ | $C^{(5)}(x)$ |
|---|---|---|---|---|---|
| 2 | 10 | 0 | 000 | 010 | 010 |
| 3 | 11 | 1 | 001 | 10 | 01 |
| 4 | 100 | 10 | 010 | 00 | 00 |
| 5 | 101 | 11 | 011 | 11 | 11 |
| 6 | 110 | 100 | 100 | 011 | 110 |

Reminder: We ultimately want to encode & decode a sequence of symbols, not just a single one (in as few bits as possible).