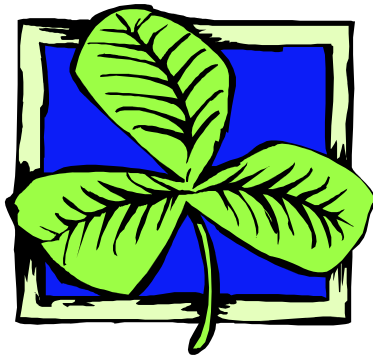


DIMACS Wkshp on Network Info. Theory  
St. Patrick's Day 2003

# Network Information Theory — Some Tentative Definitions



**James L. Massey**

Prof.-em. ETH Zurich  
Adjunct Prof., Lund Univ.  
Trondhjemsgade 3, 2TH  
DK-2100 Copenhagen East



**JamesMassey@compuserve.com**

Reminder of Gallager's notation:

$\mathbf{X}^N = X_1, X_2, \dots, X_N$  (sequence of  $N$  random variables)

$\mathbf{x}^N = x_1, x_2, \dots, x_N$  (a realization of  $\mathbf{X}^N$ )

$P(\mathbf{x}^N) = P_{\mathbf{X}^N}(\mathbf{x}^N) = \Pr(\mathbf{X}^N = \mathbf{x}^N)$

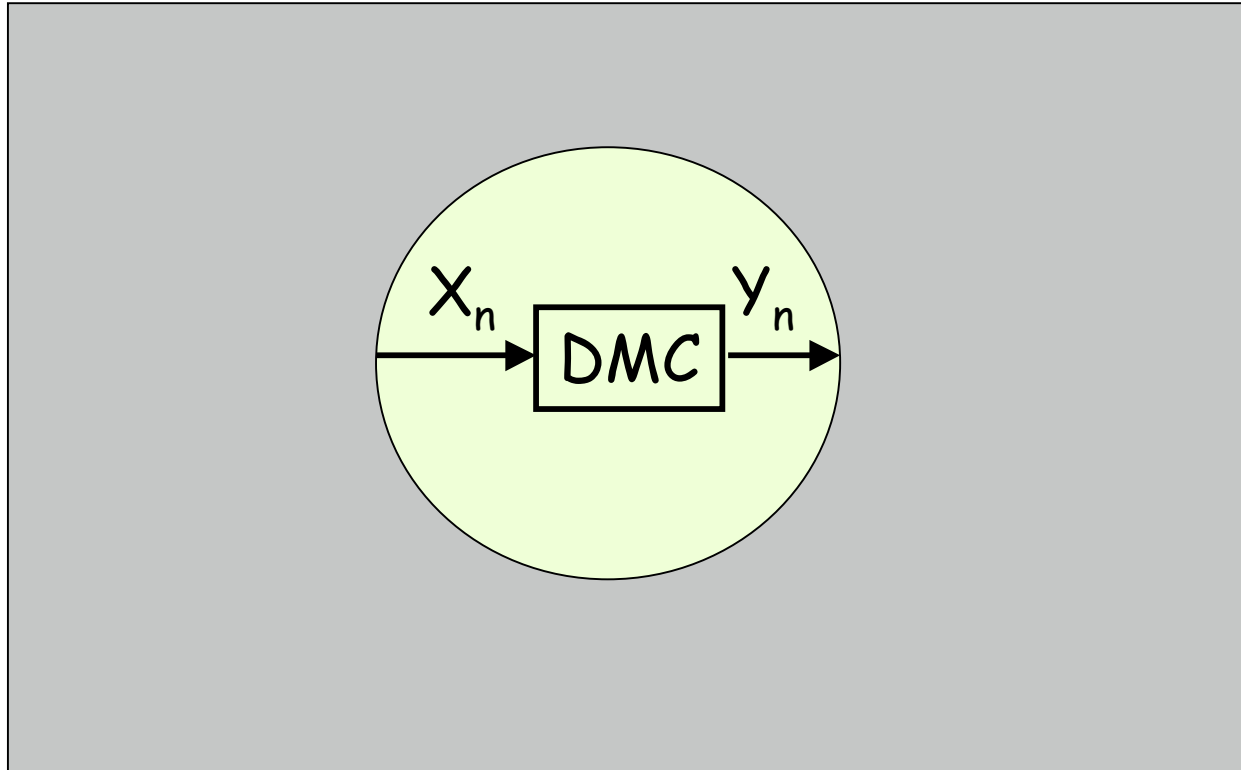
$P(\mathbf{y}^N | \mathbf{x}^N) = P_{\mathbf{Y}^N | \mathbf{X}^N}(\mathbf{y}^N | \mathbf{x}^N) = \Pr(\mathbf{Y}^N = \mathbf{y}^N | \mathbf{X}^N = \mathbf{x}^N)$

Some additional notation:

\* = "concatenation of sequences"

$\Rightarrow \mathbf{0} * \mathbf{y}^{N-1} = \mathbf{0}, y_1, y_2, \dots, y_{N-1}$





What is the probability law of a Discrete Memoryless Channel (DMC)?

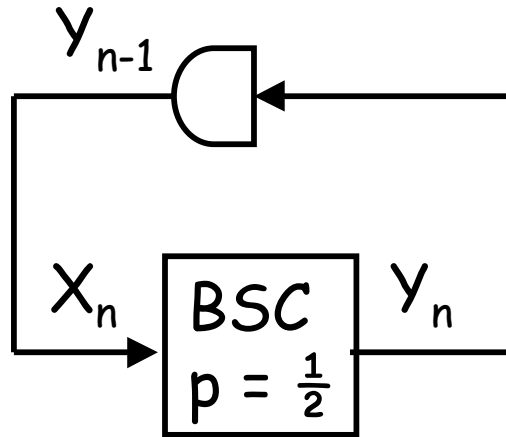
$$P(\mathbf{y}^N | \mathbf{x}^N) = \prod_{n=1}^N P_{y|X}(y_n | x_n), \text{ all } n \geq 1 \quad ???$$

Does it matter what is in the gray area?

The **natural and least restrictive assumption** to make about the gray area is that it contains **no negative delays** and **no closed paths through boxes for which the path delay is zero**.

This is the kind of assumption made in the study of discrete-time systems.





Let  $Y_0 = 0$  be the initial condition in the delay.

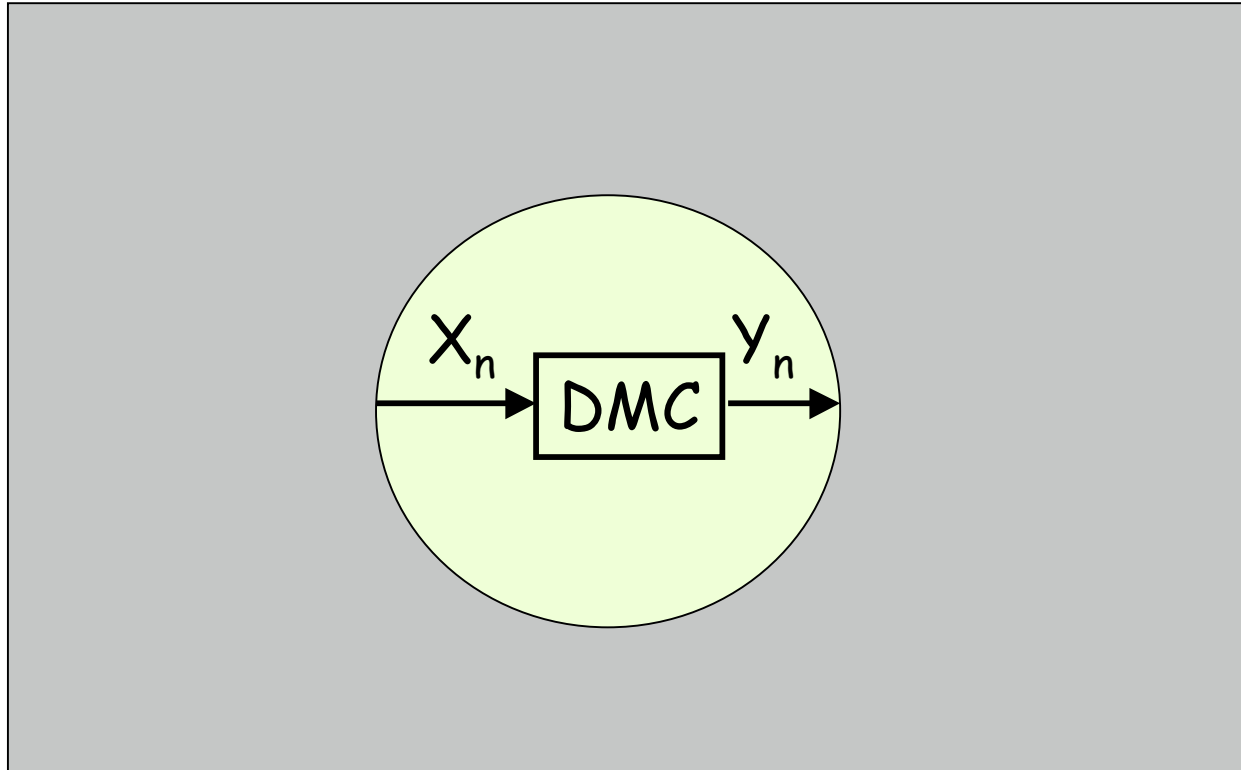
N.B.:  $X_n = Y_{n-1}$ , all  $n \geq 1$

$$\Rightarrow P(Y_1 Y_2 = 1 \ 0 \mid X_1 X_2 = 0 \ 0) = 0$$

$$\text{But } P_{Y|X}(1|0)P_{Y|X}(0|0) = 1/4$$

**What is wrong** with our definition of a DMC?





Here is the “correct” probability law of a DMC !

$$P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}) = P_{Y|X}(y_n | x_n), \text{ all } n \geq 1$$

This means that  $X_n$  is a **sufficient statistic** for reasoning about  $Y_n$  given the **observation**  $(\mathbf{X}^n, \mathbf{Y}^{n-1})$ .

In his book [R. Ash, *Information Theory*. New York: Wiley Interscience, 1965], Bob Ash correctly gave the probability law for the DMC but “spoiled” his definition by further requiring the “**causality condition**” that, for  $1 \leq n \leq N$ ,

$$P(y_n | \mathbf{x}^N \mathbf{y}^{n-1}) = P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}).$$



Adding this condition gives

$$\begin{aligned} P(\mathbf{y}^N | \mathbf{x}^N) &= \prod_{n=1}^N P(y_n | \mathbf{x}^N \mathbf{y}^{n-1}) = \prod_{n=1}^N P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}) \\ &= \prod_{n=1}^N P(y_n | \mathbf{x}_n). \end{aligned}$$

Ash’s “**causality condition**” has nothing to do with causality! What it does do is to **prohibit feedback!**

Note that Ash's "causality condition"

$$P(y_n | \mathbf{x}^N \mathbf{y}^{n-1}) = P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}) \text{ for } 1 \leq n \leq N$$

can equivalently be written as

$$H(Y_n | \mathbf{X}^N \mathbf{Y}^{n-1}) = H(Y_n | \mathbf{X}^n \mathbf{Y}^{n-1}) \text{ for } 1 \leq n \leq N.$$

It is often convenient to work with uncertainties (discrete entropies) rather than with the probability distributions themselves.

A more natural definition for prohibiting feedback is to say that **a channel is used without feedback** if

$$P(x_n | \mathbf{x}^{n-1} \mathbf{y}^{n-1}) = P(x_n | \mathbf{x}^{n-1}) \text{ for all } n \geq 1.$$



To show that these two conditions for prohibiting the use of feedback are equivalent, it suffices to show that both give the same  $P(\mathbf{x}^N \mathbf{y}^N)$ .

Using the "natural condition", we have

$$\begin{aligned} P(x_n y_n | \mathbf{x}^{n-1} \mathbf{y}^{n-1}) &= P(x_n | \mathbf{x}^{n-1} \mathbf{y}^{n-1}) P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}) \\ &= P(x_n | \mathbf{x}^{n-1}) P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}) \end{aligned}$$

$$\Rightarrow P(\mathbf{x}^N \mathbf{y}^N) = P(\mathbf{x}^N) \prod_{n=1}^N P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}).$$

Using Ash's condition, we have

$$\begin{aligned} P(\mathbf{x}^N \mathbf{y}^N) &= P(\mathbf{x}^N) P(\mathbf{y}^N | \mathbf{x}^N) = P(\mathbf{x}^N) \prod_{n=1}^N P(y_n | \mathbf{x}^N \mathbf{y}^{n-1}) \\ &= P(\mathbf{x}^N) \prod_{n=1}^N P(y_n | \mathbf{x}^n \mathbf{y}^{n-1}). \end{aligned}$$

Unlike causality, probabilistic dependence has no direction. Whether A causes B or B causes A, A and B will be statistically dependent.

Essentially this is why  $I(\mathbf{X}^N; \mathbf{Y}^N) = I(\mathbf{Y}^N; \mathbf{X}^N)$

or, equivalently,

$$H(\mathbf{Y}^N) - H(\mathbf{Y}^N | \mathbf{X}^N) = H(\mathbf{X}^N) - H(\mathbf{X}^N | \mathbf{Y}^N).$$



My 1990 definition of **directed information**:

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) = \sum_{n=1}^N I(\mathbf{X}^n; Y_n | \mathbf{Y}^{n-1}).$$

(By ignoring information that  $Y_n$  may be giving about future  $X$  digits, we are considering only the information flowing from  $X$  digits to  $Y$  digits.)

In terms of uncertainties (discrete entropies)

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) = \sum_{n=1}^N [H(Y_n | \mathbf{Y}^{n-1}) - H(Y_n | \mathbf{X}^n \mathbf{Y}^{n-1})].$$

We recall that

$$I(\mathbf{X}^N; \mathbf{Y}^N) = \sum_{n=1}^N [H(Y_n | \mathbf{Y}^{n-1}) - H(Y_n | \mathbf{X}^N \mathbf{Y}^{n-1})].$$

Equivalently, we can write

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) = \sum_{n=1}^N [H(\mathbf{X}^n | \mathbf{Y}^{n-1}) - H(\mathbf{X}^n | \mathbf{Y}^n)].$$

Because  $H(Y_n | \mathbf{X}^N \mathbf{Y}^{n-1}) \leq H(Y_n | \mathbf{X}^n \mathbf{Y}^{n-1})$ ,  
it follows that

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) \leq I(\mathbf{X}^N ; \mathbf{Y}^N)$$

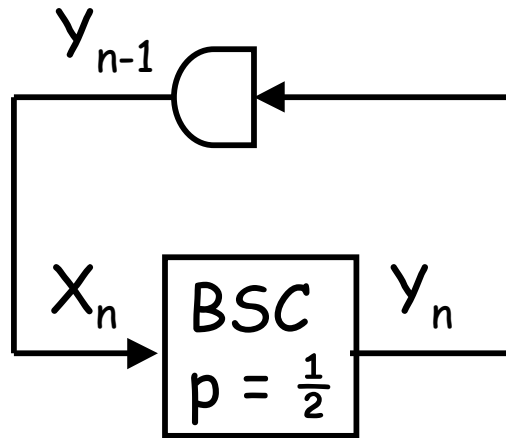
with equality if and only if there is no feedback, i.e., if and only if (natural condition)  $P(x_n | \mathbf{x}^{n-1} \mathbf{y}^{n-1}) = P(x_n | \mathbf{x}^{n-1})$  [or, equivalently, (Ash's condition)  $P(y_n | \mathbf{x}^N \mathbf{y}^{n-1}) = P(y_n | \mathbf{x}^n \mathbf{y}^{n-1})$ ] for  $1 \leq n \leq N$ .

If  $\mathbf{X}^N$  and  $\mathbf{Y}^N$  are the input and output sequences of a DMC, then

$$\begin{aligned}
 I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) &= \sum_{n=1}^N [H(Y_n | \mathbf{Y}^{n-1}) - H(Y_n | \mathbf{X}^n \mathbf{Y}^{n-1})] \\
 &= \sum_{n=1}^N [H(Y_n | \mathbf{Y}^{n-1}) - H(Y_n | X_n)] \\
 &\leq \sum_{n=1}^N [H(Y_n) - H(Y_n | X_n)] = \\
 &= \sum_{n=1}^N I(X_n; Y_n)
 \end{aligned}$$



with equality if and only if  $Y_1, Y_2, \dots, Y_N$  are statistically independent.



Let  $Y_0 = 0$  be the initial contents of the delay.

N.B.:  $X_n = Y_{n-1}$ , all  $n \geq 1$   
so that  $X_1 = 0$

$$H(Y^N) = N \text{ bits}$$

$$H(Y^N | X^N) = H(Y_N | X^N) = 1 \text{ bit}$$

$$\Rightarrow I(X^N; Y^N) = N - 1 \text{ bits}$$

$$I(X_n; Y_n) = 0 \text{ bits, all } n$$

$$\Rightarrow I(X^N \rightarrow Y^N) = 0 \text{ bits}$$

$$H(X_n | X^{n-1}) = 1 \text{ bit, all } n \geq 2$$

$$H(X_1) = 0 \text{ bits}$$

$$H(X_n | O^* Y^{n-1} X^{n-1}) = 0, n > 1$$

$$\Rightarrow I(O^* Y^{N-1} \rightarrow X^N) = N - 1 \text{ bits}$$

Note that  $I(X^N \rightarrow Y^N) + I(O^* Y^{N-1} \rightarrow X^N) = I(X^N; Y^N)$ .

## Conservation Law for Directed Information

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) + I(\mathbf{O}^* \mathbf{Y}^{N-1} \rightarrow \mathbf{X}^N) = I(\mathbf{X}^N; \mathbf{Y}^N).$$



Proof by induction:

$$I(\mathbf{X}^1 \rightarrow \mathbf{Y}^1) = I(X_1; Y_1) \text{ and } I(\mathbf{O} \rightarrow \mathbf{X}^1) = 0.$$

Note that

$$I(\mathbf{X}^{n+1} \rightarrow \mathbf{Y}^{n+1}) = I(\mathbf{X}^n \rightarrow \mathbf{Y}^n) + I(\mathbf{X}^{n+1}; Y_{n+1} | \mathbf{Y}^n)$$

and similarly that

$$\begin{aligned} I(\mathbf{O}^* \mathbf{Y}^n \rightarrow \mathbf{X}^{n+1}) &= I(\mathbf{O}^* \mathbf{Y}^{n-1} \rightarrow \mathbf{X}^n) + I(\mathbf{O}^* \mathbf{Y}^n; X_{n+1} | \mathbf{X}^n) \\ &= I(\mathbf{O}^* \mathbf{Y}^{n-1} \rightarrow \mathbf{X}^n) + I(\mathbf{Y}^n; X_{n+1} | \mathbf{X}^n). \end{aligned}$$

Thus, by the induction hypothesis,

$$\begin{aligned} I(\mathbf{X}^{n+1} \rightarrow \mathbf{Y}^{n+1}) + I(\mathbf{O}^* \mathbf{Y}^n \rightarrow \mathbf{X}^{n+1}) &= I(\mathbf{X}^n; \mathbf{Y}^n) + I(\mathbf{X}^{n+1}; Y_{n+1} | \mathbf{Y}^n) \\ &\quad + I(\mathbf{Y}^n; X_{n+1} | \mathbf{X}^n) \\ &= I(\mathbf{X}^{n+1}; \mathbf{Y}^n) + I(\mathbf{X}^{n+1}; Y_{n+1} | \mathbf{Y}^n) \\ &= I(\mathbf{X}^{n+1}; \mathbf{Y}^{n+1}). \end{aligned}$$

Gerhard Kramer's definition of **causal conditioning**:

$$H(\mathbf{Y}^N \parallel \mathbf{X}^N) = \sum_{n=1}^N H(Y_n \mid \mathbf{X}^n \mathbf{y}^{n-1})$$



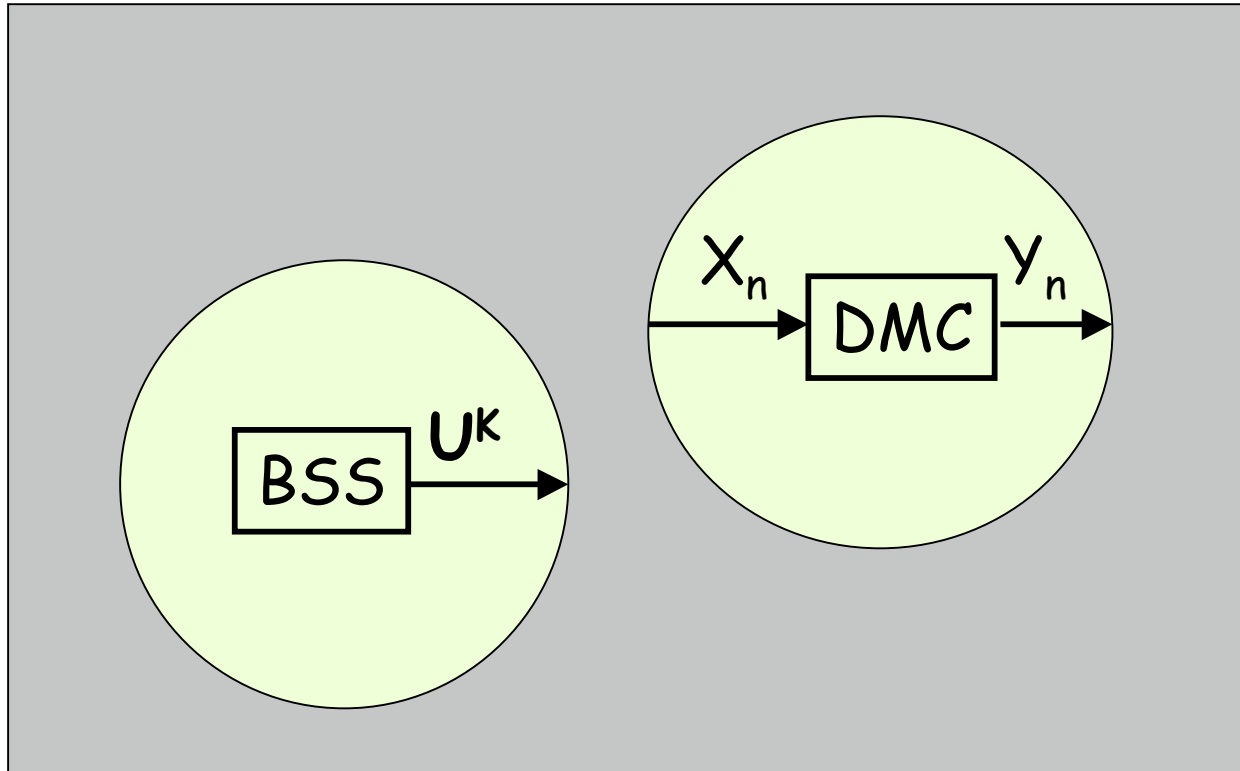
In terms of causal conditioning, one can write

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) = H(\mathbf{Y}^N) - H(\mathbf{Y}^N \parallel \mathbf{X}^N)$$

Gerhard also defined causally conditioned directed information in the following manner:

$$I(\mathbf{X}^N \rightarrow \mathbf{Y}^N \parallel \mathbf{Z}^N) = H(\mathbf{Y}^N \parallel \mathbf{Z}^N) - H(\mathbf{Y}^N \parallel \mathbf{X}^N \mathbf{Z}^N)$$





We want to consider a synchronized network in which all devices are governed by the same notion of time.



What we mean by **synchronization**:

$X_1, X_2, X_3, X_4, X_5, \dots, X_N$   
 $Y_1, Y_2, Y_3, Y_4, Y_5, \dots, Y_N$   
 $Z_1, Z_2, Z_3, Z_4, Z_5, \dots, Z_N$

Three arbitrary  
**clocked sequences**

The meaning is that the  $n^{\text{th}}$  variable in each sequence, i.e.,  $X_n$ ,  $Y_n$ , and  $Z_n$ , take on their values at the same time instant. Moreover, the  $n^{\text{th}}$  variable takes on its value before the  $(n+1)^{\text{st}}$  variable. Note also that  $Y_{n-1}$  is the  $n^{\text{th}}$  variable in the concatenated sequence  $0 * Y^{N-1}$ .

**Which kinds of sequences should be clocked?**

**Channels:** Input and output sequences are clocked.

**Delays:** Input and output sequences are clocked.

**Sources:** **Not clocked** - output present at creation!

**Channel Encoders:** Only output sequence and input feedback sequence (if present) are clocked.

**Channel Decoders:** Only input sequence and output feedback sequence (if present) are clocked.

**Source Encoders:** **Not clocked** - input and output present at creation!

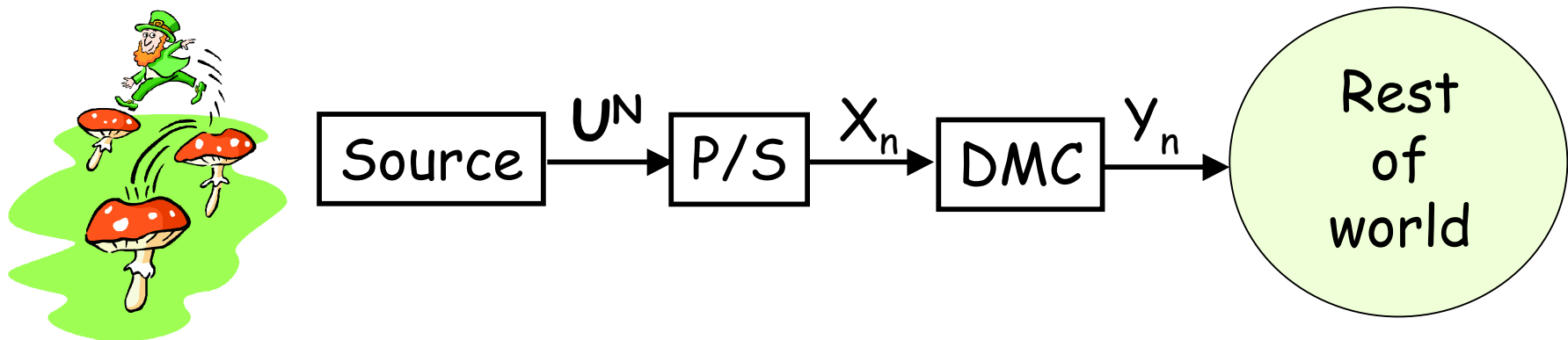
**Source Reconstructors:** **Not clocked** - input and output present at end of world!



Clocked inputs to network components can come only from clocked outputs of other network components.

⇒ you cannot connect a source directly to a channel!  
You must use a channel encoder.

Michael Gastpar's direct transmission of a source over a channel has to be understood as using a parallel-to-series converter as the channel encoder.



The **causality assumption** for synchronized networks:

For any channel (whether or not memoryless), the input-output sequences  $(X^n, Y^{n-1})$  are a **sufficient statistic** for reasoning about the output  $Y_n$  given the **observation**  $(X^n, Y^{n-1}, W)$ , where  $W$  is any random quantity composed of inputs and outputs of other network elements at time  $n$  or earlier,  $1 \leq n \leq N$ .

Thus, for any **source output**  $U^K$ ,

$$H(Y_n | X^n Y^{n-1} U^K) = H(Y_n | X^n Y^{n-1}).$$

$$\begin{aligned} \Rightarrow H(Y^N | U^K) &= \sum_{n=1}^N H(Y_n | U^K Y^{n-1}) \geq \sum_{n=1}^N H(Y_n | U^K X^n Y^{n-1}) = \\ &= \sum_{n=1}^N H(Y_n | X^n Y^{n-1}) = H(Y^N || X^N) \end{aligned}$$

Thus,  **$I(U^K ; Y^N) \leq I(X^N \rightarrow Y^N)$** .

An immediate consequence of the fact that

$$I(\mathbf{U}^K ; \mathbf{Y}^N) \leq I(\mathbf{X}^N \rightarrow \mathbf{Y}^N) \leq \sum_{n=1}^N I(X_n ; Y_n) \leq N C_{\text{DMC}}$$

is that **feedback does not increase the capacity of a discrete memoryless channel (DMC).**

How can one logically prove this result if one uses the “usual definition” of a DMC, which is given on slide 2 and which in fact prohibits the use of feedback?

How can we reason correctly about complicated networks of sources, channel, encoders, decoders and delays without some sort of careful statement of our assumptions similar to that given here today?



## Some references:

J. L. Massey, "Causality, Feedback and Directed Information," pp. 303-305 in *Proc. 1990 Int. Symp. on Info. Th. & its Appls.*, Hawaii, USA, Nov. 27-30, 1990.

H. Marko, "The Bidirectional Communication Theory – A Generalization of Information Theory", *IEEE Trans. Commun.*, vol. COM-21, pp. 1345-1351, Dec. 1973.

G. Kramer, *Directed Information for Channels with Feedback*, ETH Series in Inform. Proc., vol. 11. Konstanz: Hartung--Gorre, 1998.

G. Kramer, "Capacity Results for the Discrete Memoryless Network," *IEEE Trans. Inform. Th.*, vol. IT-49, pp. 4-21, Jan. 2003.