

# How is Digital Sequence Information (DSI) produced ?

Dr Wael Houssen

EPSRC Fellow

Institute of Medical Sciences,  
University of Aberdeen.

<https://www.abdn.ac.uk/people/w.houssen>





CBD



**Convention on  
Biological Diversity**

Distr.  
GENERAL

CBD/DSI/AHTEG/2020/1/3  
29 January 2020

ENGLISH ONLY

AD HOC TECHNICAL EXPERT GROUP ON  
DIGITAL SEQUENCE INFORMATION ON  
GENETIC RESOURCES

Montreal, Canada, 17-20 March 2020

*Annex*

**Digital Sequence Information on Genetic Resources: Concept, Scope and  
Current Use**

*Wael Housen, Rodrigo Sara, Marcel Jaspars*

**Links:**

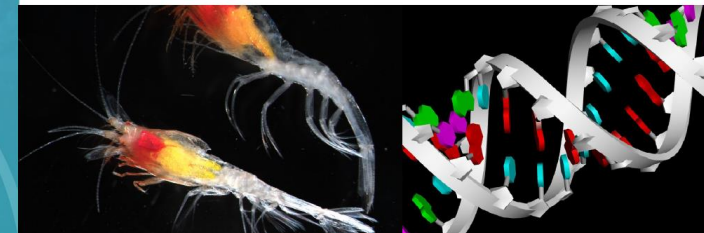
- <https://www.cbd.int/doc/c/fef9/2f90/70f037ccc5da885dfb293e88/dsi-ahteg-2020-01-03-en.pdf>
- <https://www.dosi-project.org/wp-content/uploads/070-DSI-Policy-brief-V4-WEB.pdf>

# POLICY BRIEF

MARCH 2020



## Digital Sequence Information – Clarifying Concepts

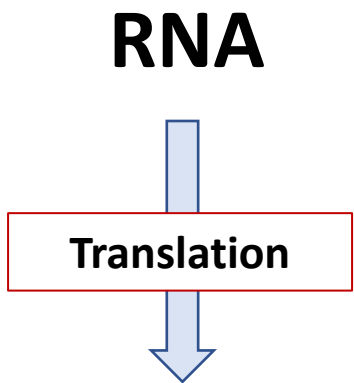
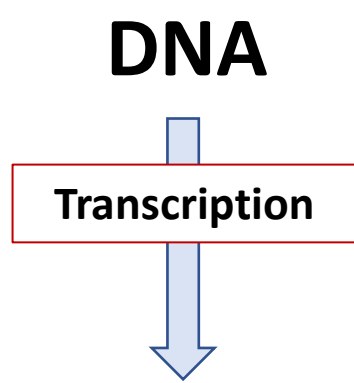


This Policy Brief is Based on:

“Digital Sequence Information on Genetic Resources: Concept, Scope and Current Use”. Authors:

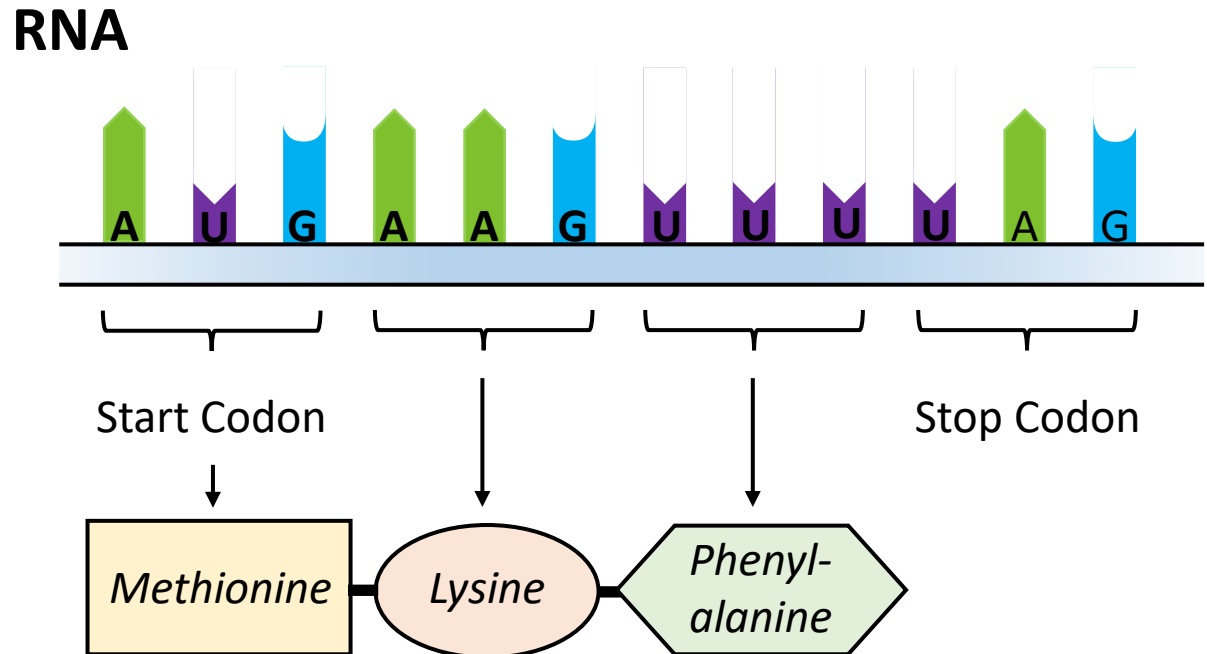
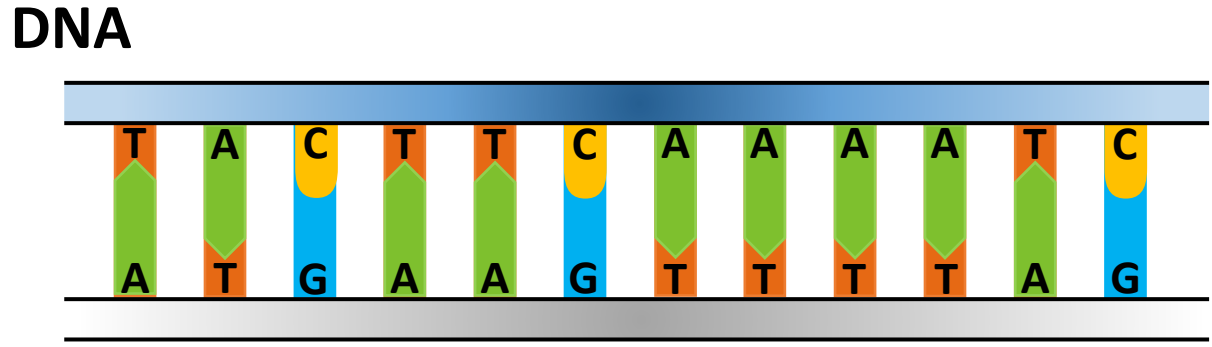
*Wael Housen, Rodrigo Sara, Marcel Jaspars*

# The Central Dogma of Molecular Biology

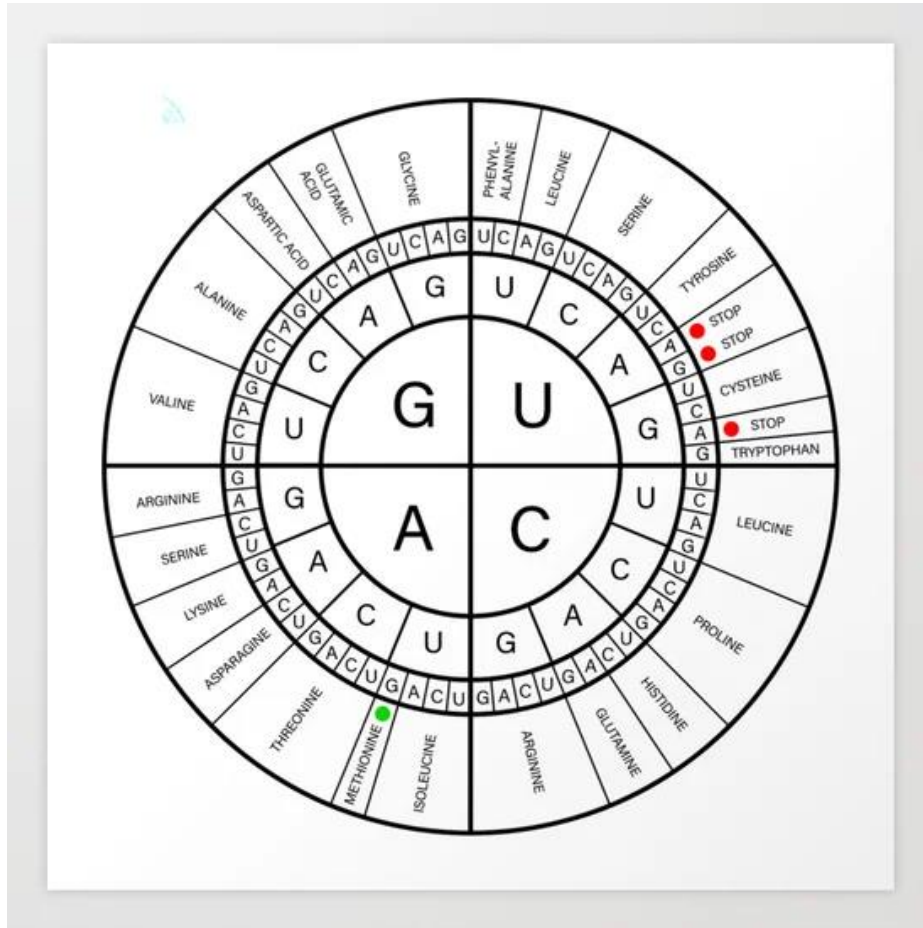


**Protein**

**Metabolites**



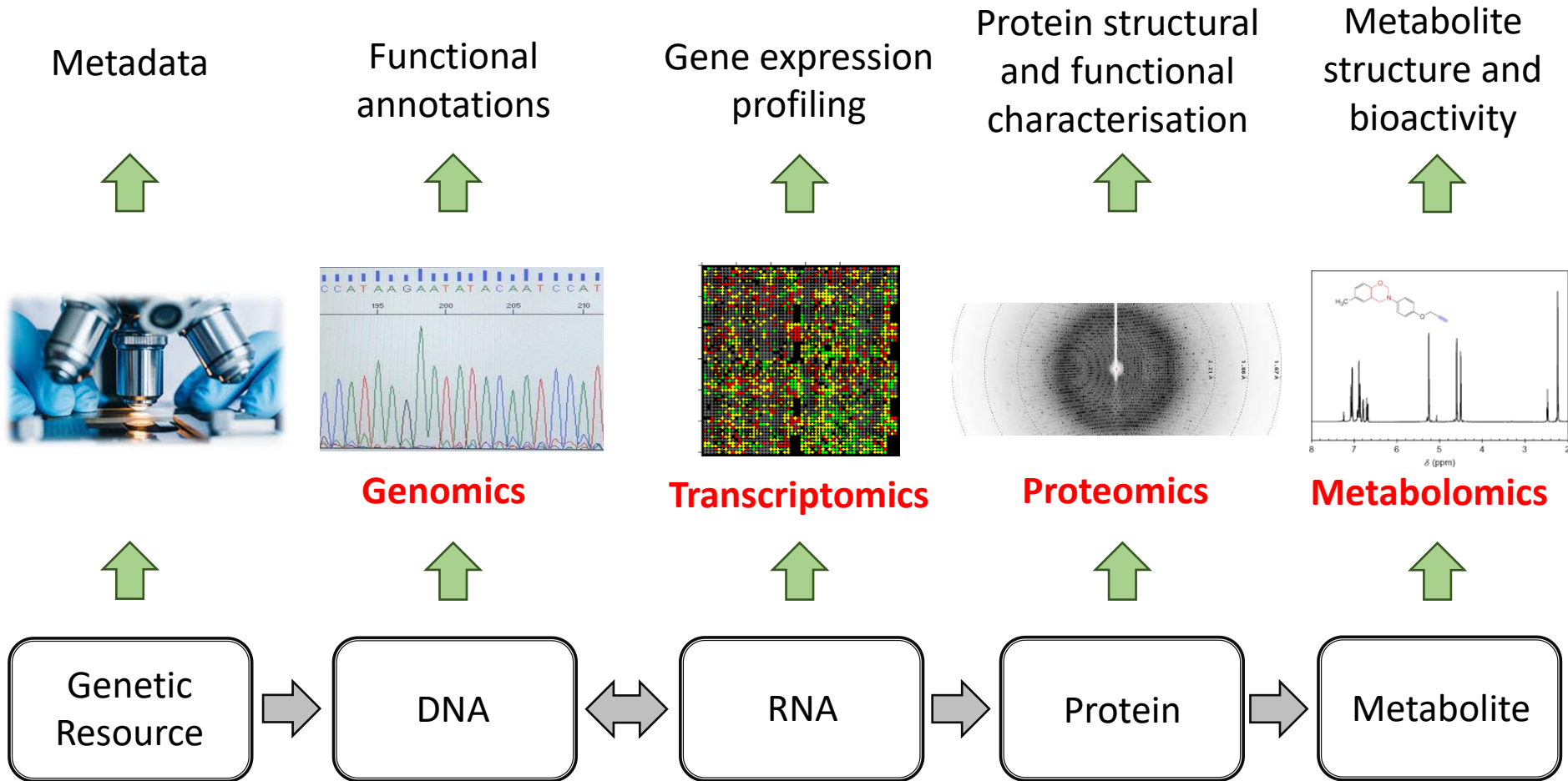
# The Codon wheel



- Sixty-four codons encode for the twenty natural amino acids.
- ‘AUG’ encodes the amino acid methionine, and activates the ribosome to start the process of making a protein and is thus known as “start codon”.
- “Stop codons” signal the termination of translation into proteins.
- Genetic code is degenerate is that many amino acids can be encoded by different codons.
- Each organism has a ‘preference’ for the use of a particular codon for a specific amino acid. This is called “Codon bias”.

# Origin of Digital Sequence Information

## Applications



# Raw sequence vs. value-added DSI

---

DNA sequence alone does not provide information on:

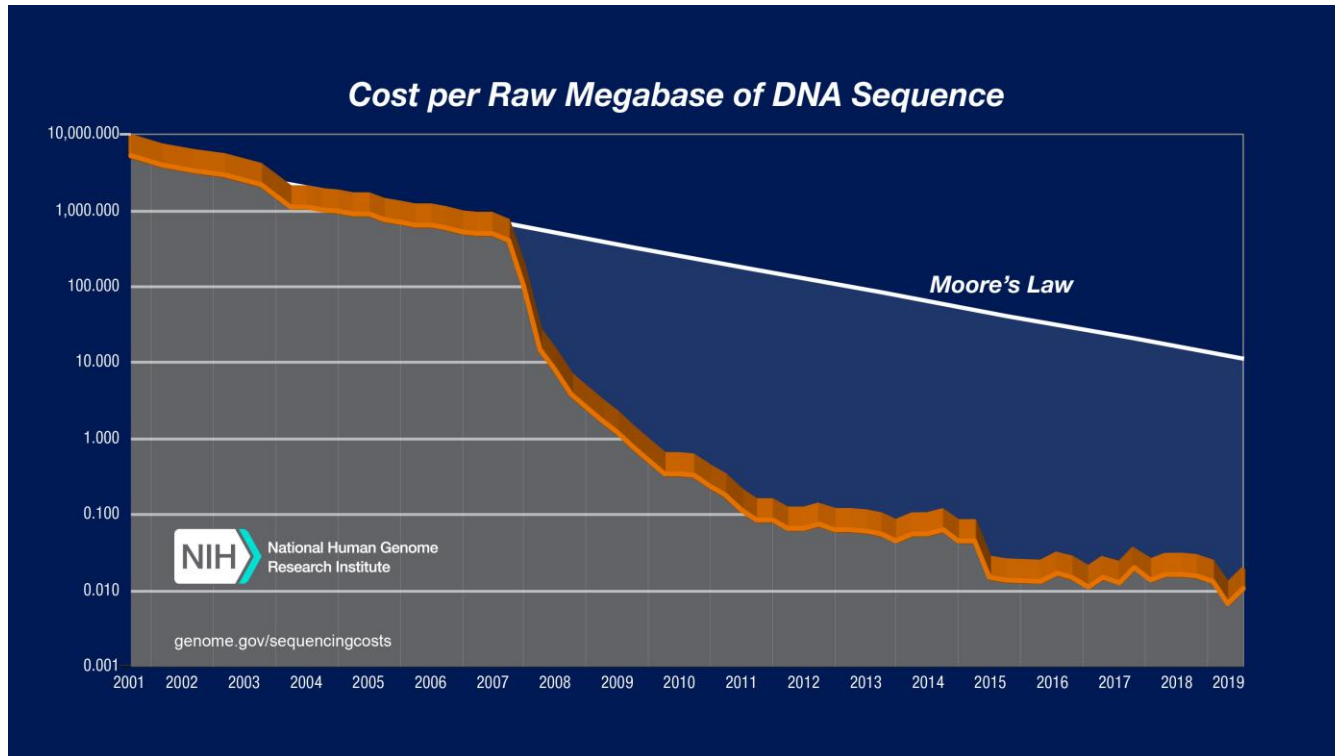
- The function of each gene.
- The level of expression of each gene.
- How the encoded protein sequence folds. Misfolded proteins are known to underlie serious diseases.
- The posttranslational modifications in proteins e.g. protein glycosylation.

# Next Generation Sequencing Technologies

Platform / Instrument	Throughput range (Gb)	Read length (bp)	Strength	Weakness
<b>Sanger sequencing</b>				
<b>ABI 3500/3730</b>	0.0003	Up to 1kb	Read accuracy and length	Cost and throughput
<b>Illumina</b>				
<b>MiniSeq</b>	1.7–7.5	1×75 to 2×150	Low initial investment	Run and read length
<b>MiSeq</b>	0.3–15	1×36 to 2×300	Read length, scalability	Run length
<b>NextSeq</b>	10–120	1×75 to 2×150	Throughput	Run and read length
<b>HiSeq (2500)</b>	10–1000	1×50 to 2×250	Read accuracy, throughput, low per sample cost	High initial investment, run length
<b>NovaSeq 5000/6000</b>	2000–6000	2×50 to 2×150	Read accuracy, throughput Low per sample cost	High initial investment, run and read length
<b>Ion Torrent</b>				
<b>PGM</b>	0.08–2	Up to 400	Read length, speed	Throughput, homopolymers
<b>S5</b>	0.6–15	Up to 400	Read length, speed, scalability	Homopolymers
<b>Proton</b>	10–15	Up to 200	Speed, throughput	Homopolymers
<b>Pacific BioSciences</b>				
<b>PacBio RSII</b>	0.5–1	Up to 60 kb (Average 10 kb, N50 20 kb)	Read length, speed	High error rate and initial investment, low throughput
<b>Sequel</b>	5–10	Up to 60 kb (Average 10 kb, N50 20 k)	Read length, speed	High error rate
<b>Oxford Nanopore</b>				
<b>MinION</b>	0.1–1	Up to 100 kb	Read length, portability	High error rate, run length, low throughput

# Sequencing costs

---



- Sequencing is getting **faster – cheaper**
- Number of sequences continues to increase.



# What is being sequenced?

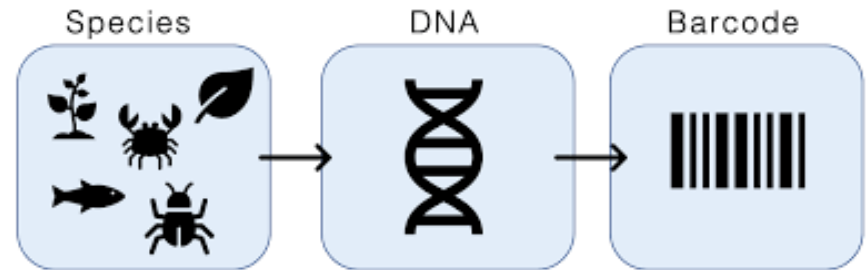
---

- Whole genome.
- Metagenome.
- 16S/18S/ITS/CO1 amplicon for molecular identification of microorganisms.
- Environmental DNA.

# Examples of the applications of DSI

---

- Taxonomical identification of organisms: DNA barcoding



- Biodiversity conservation
- Agriculture and food security: Genetic markers to assist selective breeding; Development of GMO
- Drug discovery: Identification of new drug targets
- Medicine: Detection of infectious pathogens
- SynBio and biotechnology: Recombinant production of proteins e.g. insulin and enzymes for laundry detergents.