

Uses of DNA Sequence Data

Marcel Jaspars

Marine Biodiscovery Centre, Department of Chemistry

University of Aberdeen, Scotland, UK

m.jaspars@abdn.ac.uk

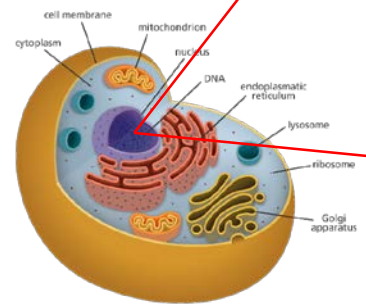
Outline

- Production of DSI
 - Genetic resources
 - The central dogma of molecular biology
 - The Origin of DSI
- The rise of 'omics'
- What is being sequenced and why
 - Example – fish identification
- DSI usage
 - Non-linearity of DSI usage
 - Example - enzyme modification to make pharmaceuticals
 - Synthetic biology to make pharmaceuticals
- Conclusions

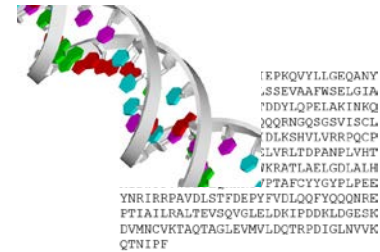
Production of DSI - Genetic Resources



Biological resource

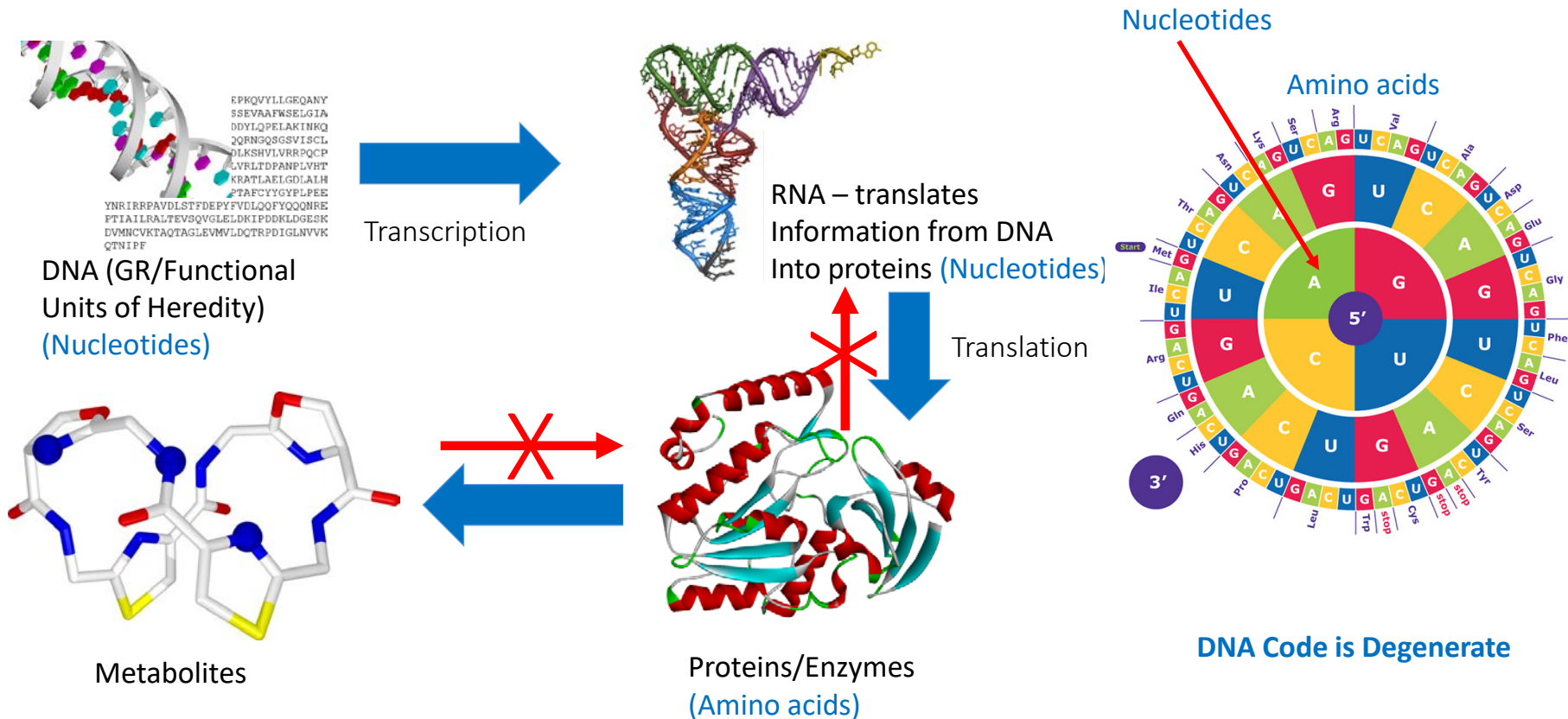


Cells (one or more)

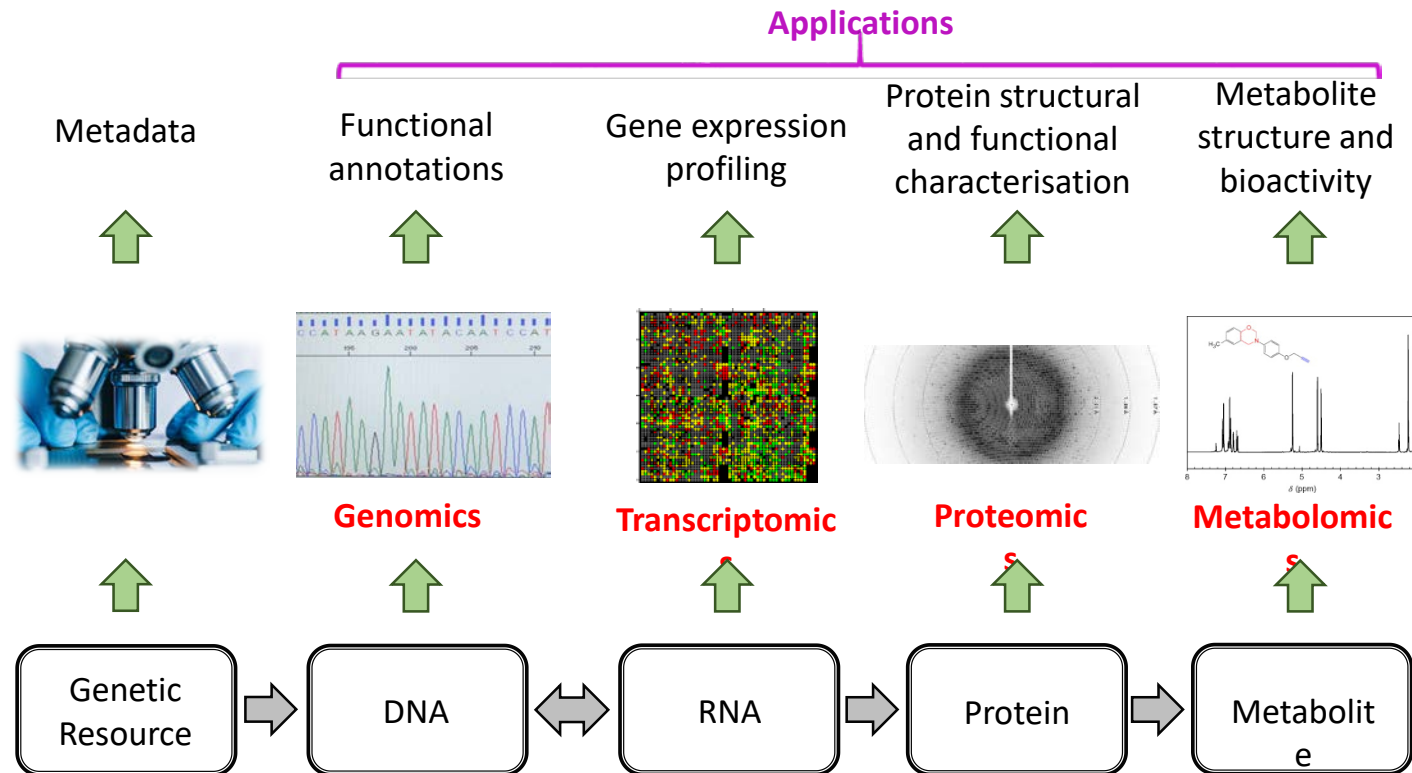


DNA (GR/Functional Units of Heredity)

Production of DSI - The Central Dogma of Molecular Biology



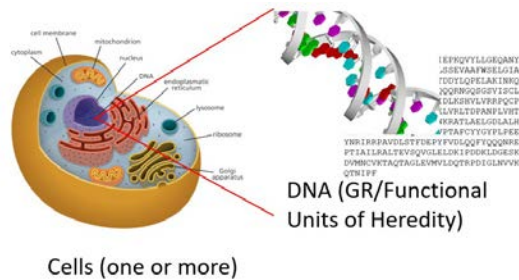
The Origin of Digital Sequence Information



The Rise of 'Omics'

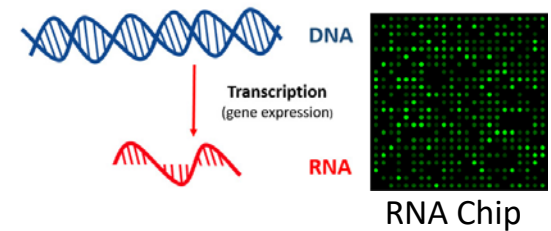
Genomics

The genome is the complete set of genes or genetic material present in a cell or organism.



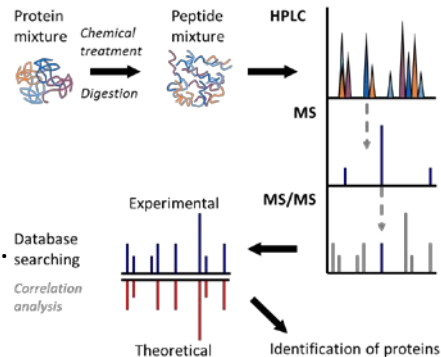
Transcriptomics

The transcriptome is a measure of which genes are expressed under any given set of conditions.



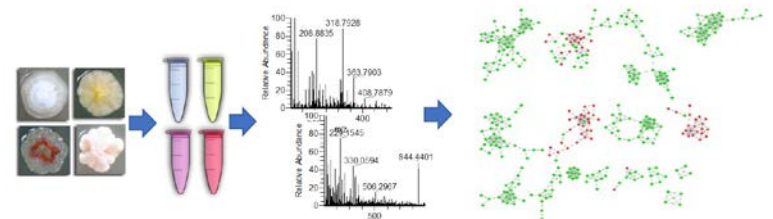
Proteomics

The proteome is the totality of all proteins produced or modified under any given set of conditions.



Metabolomics

The metabolome is the totality of small molecule metabolites produced under any given set of conditions.



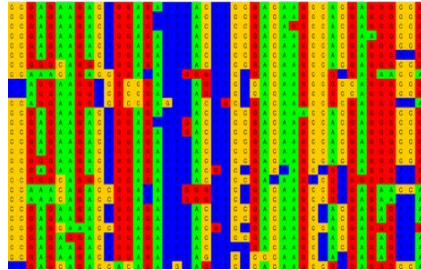
What is Being Sequenced and Why?

Ecosystems



Understand how life forms interact with each other and their environment.

Populations



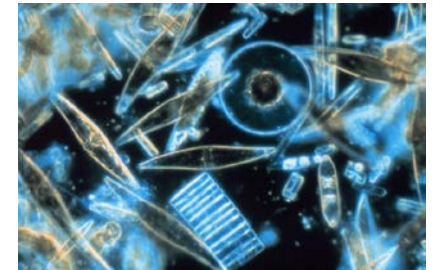
Understand variation within or between related species to learn more about their biology through comparative genomics.

Biodiversity



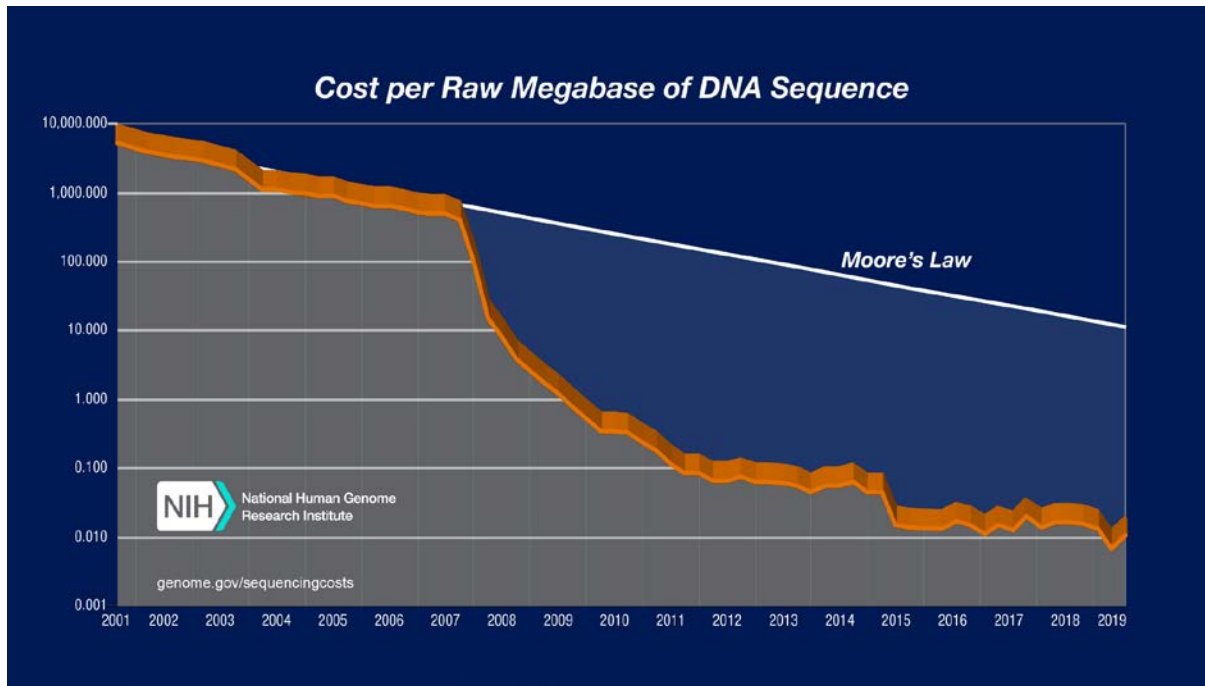
Learn about life forms and preserve the genetic data before they become extinct.

Metagenomes



Understanding the microbial communities and their environments

Sequencing Costs



Cost of sequencing continues to decline

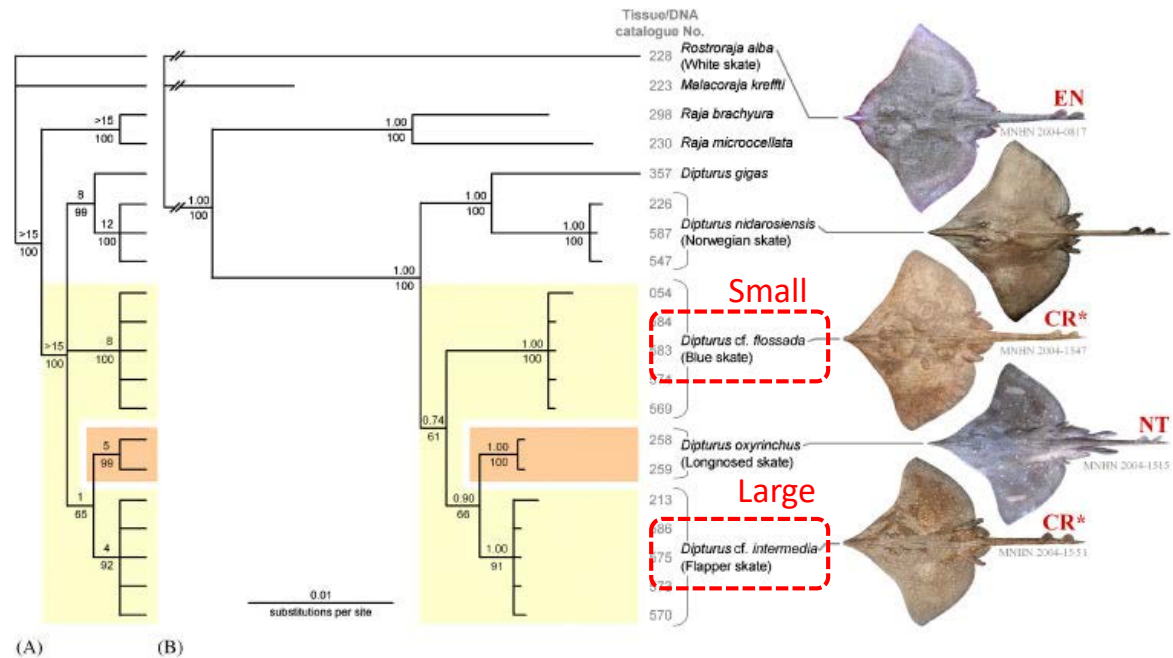
Number of sequences escalates – **data avalanche**

Accurate interpretation and annotation of sequences needs to improve

Example of Application of DNA Sequences

Mistaken Identity of Fish

- Declining Catches of Common Skate '*Dipturus batis*'
- Taxonomic discrepancies
- mtDNA gave phylogenetic tree
- Two species
 - Small one common
 - Large one rare
- Information used for setting catch limits

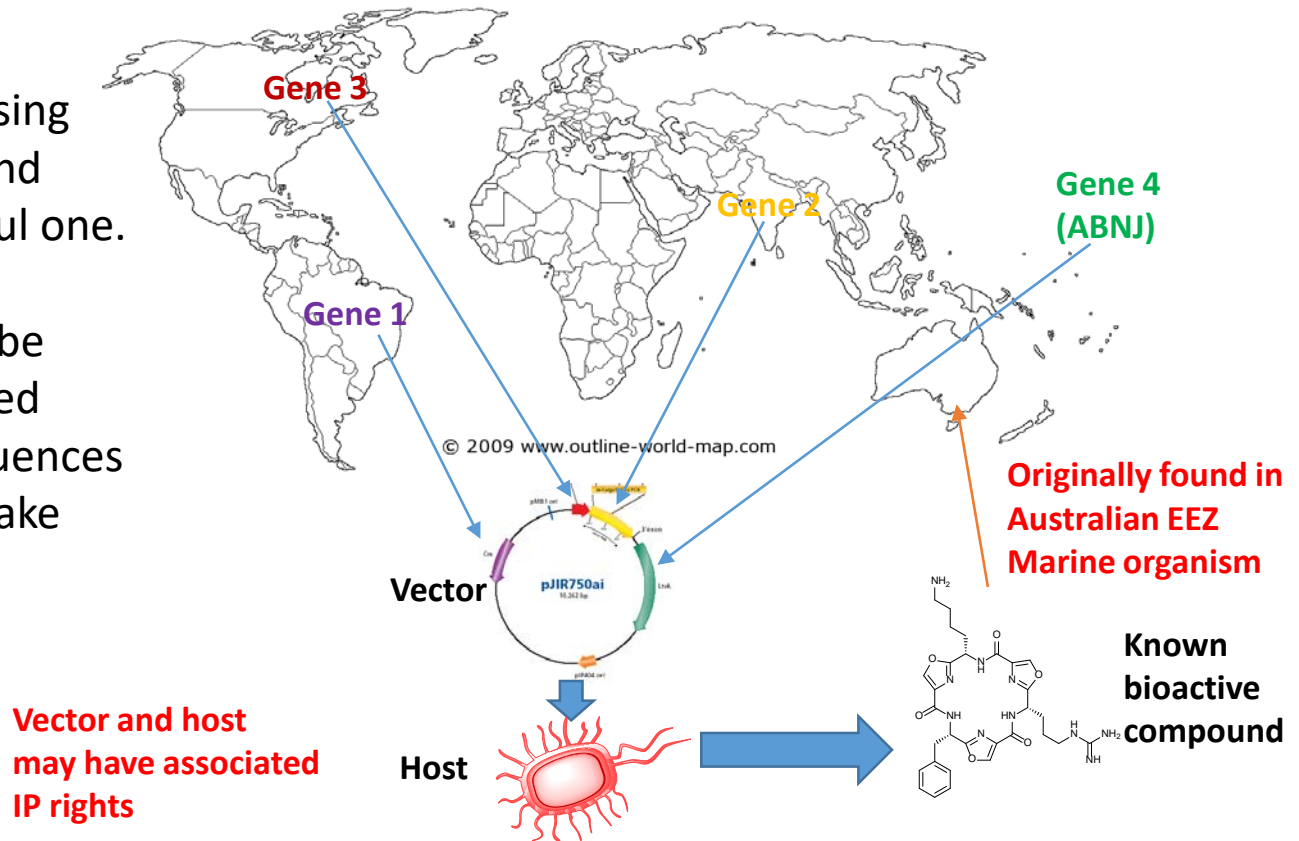


Aquatic Conserv: Mar. Freshw. Ecosyst. 20: 319–333.

DSI Usage is Not Straightforward or 'Linear'

Sequences used can be the result of analysing multiple sequences and finding the most useful one.

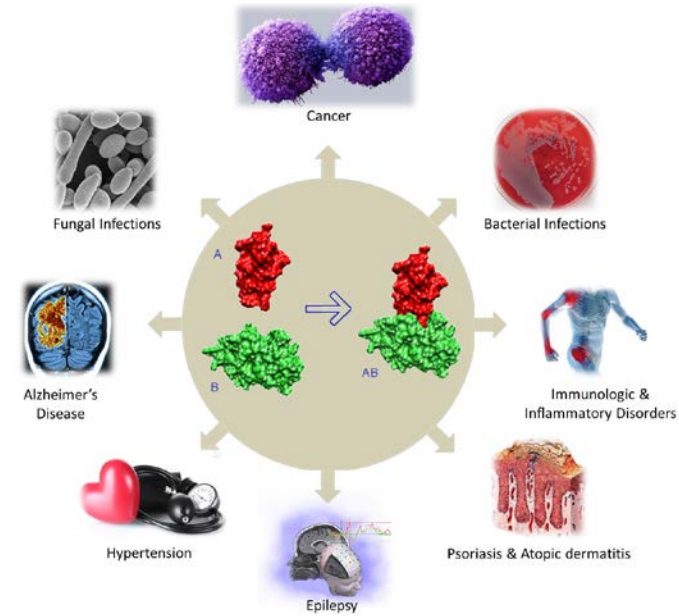
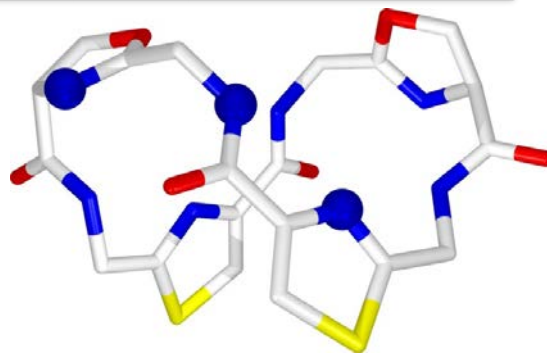
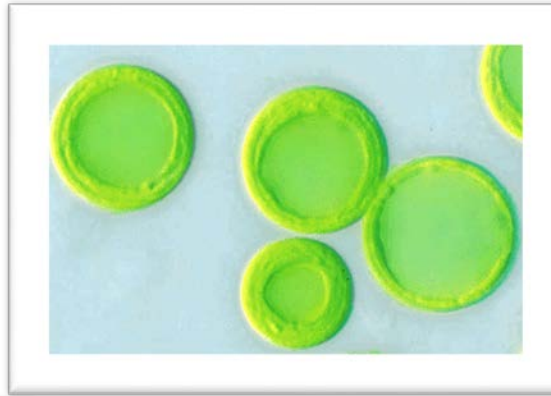
Some sequences can be designed or engineered based on natural sequences and the process can take several iterations.





Gyrocycle™ highly modified macrocyclic peptides - effective alternatives to therapeutic peptides

Taking the Lead from Nature

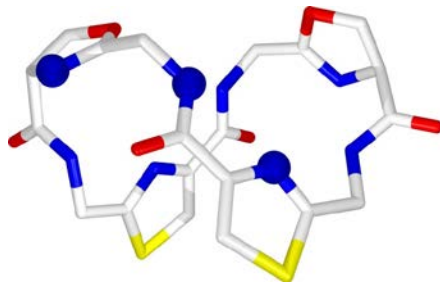


Macrocycles may be effective against a range of complex diseases

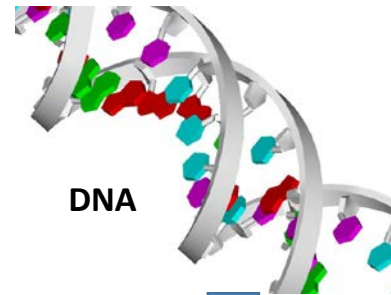


Gyrocycle™ highly modified macrocyclic peptides - effective alternatives to therapeutic peptides

Using DNA Sequence Data



Macrocycle



DNA



Protein



Using Protein Sequence Data

```

MQPTALQIKPHFHVEIIEPKQVYLLGEEQGNHALTGQLY
LSRLVEKGYLTVAPELSLEVAAFWSELGIAPSVVAEG
VSDPKAPKAPKAGDSTAQLQVVLTDDYLQPELAAINKE
HCLAQRLRGNREVEASVLQKRALQERNQKNGAVSC
NAIAPGTARFPTLAGKIFTFNQTTLELKAHPLSRPQC
RATTPQQTQVKYQHLIGPITGVVTELVRI SDPANPLVH
DSQSRASGLCEAIERYSIGIFLGDEPRKRATLAEGLDLA
PHRFAASQAIDWTPWLSLTEQKHKYVPTAICYNYLLF
RDSVALWYNRLRRPEVELSSFEPEYFLQLQQFYRSQN
IGFGAHLDPKIAILLRALTEVSQVGLLELDKVPDEKLDGE
RWSDDIYTDVMACVEMAKVAGLETVLVDQTRPDIGLNV
PLAEAE MNPTNIPF
    
```

Original Protein Sequence



Native enzyme



```

MQSTLLQIKPHFHVEIEPKQVYLLGEEQGNHALTG
VLNRLAEKGYLTEATPDLSEVAAFWTELGIAPTVA
PVQNAS MQSTPLLQIQPHFHVEIEPKQVYLLGEQANY
RLRGNR VLDRLAEKGYLTEAAPELSSEVAAFWSELGIA
TALFPT VQTSTEAGSPTALNVVLTDDYLQPELAKINKQ
QTLQKY LRGNRMQSTPLLQIQPHFHVEIEPKQVYLLGEQANY
GFCEAV VFFPT: VLDRLAEKGYLTEAAPELSSEVAAFWSELGIA
QAIDWT TVQKY: VQTSTEAGSPTALNVVLTDDYLQPELAKINKQ
WYNRLS LCEAI: LRGNREVEASVLQKQAQQQRNGQSGSVISCL
DPTTGI AIDWT: VFFPTLDGKIITFNHTVIDLKSHVLVRRPQCP
TDVMTN YNRIR: TVQKYQHLVSPITGVVTELVRLTDPANPLVHT
NPMNIP PTIAI: LCEAIERYSIGIFQGDEPWKRATLAEGLDLALH
DVMNC: AIDWTPVWSLTEQKHKYVPTAFCYYGYPLPEE
QTNIP: YNRIRRPVDLSTFDEPYFVDLQQFYQQQNRRE
PTIAILLRALTEVSQVGLLELDKIPDDKLDGESK
DVMNCVKTAQTAGLEVMLVDQTRPDIGLNVVK
QTNIPF
    
```

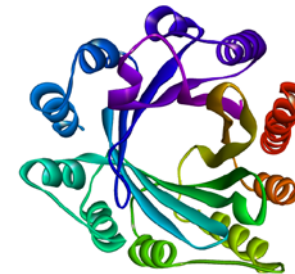
Use DNA databases to find analogues and make comparisons



```

MQPTALQIKPHFHVEIIEPKQVYLLGEEQGNHALTGQLY
LSRLVEKGYLTVAPELSLEVAAFWSELGIAPSVVAEG
VSDPKAPKAPKAGDSTAQLQVVLTDDYLQPELAAINKE
HCLAQRLRGNREVEASVLQKRALQERNQKNGAVSC
NAIAPGTARFPTLAGKIFTFNQTTLELKAHPLSRPQC
RATTPQQTQVKYQHLIGPITGVVTELVRI SDPANPLVH
DSQSRASGLCEAIERYSIGIFLGDEPRKRATLAEGLDLA
PHRFAASQAIDWTPWLSLTEQKHKYVPTAICYNYLLF
RDSVALWYNRLRRPEVELSSFEPEYFLQLQQFYRSQN
IGFGAHLDPKIAILLRALTEVSQVGLLELDKVPDEKLDGE
RWSDDIYTDVMACVEMAKVAGLETVLVDQTRPDIGLNV
PLAEAE MNPTNIPF
    
```

Choose sequence(s)
Targeted modifications
Gene Synthesis



Engineered Enzyme
Desired T range
Desired pH range
Desired specificity

Can be produced at High titres



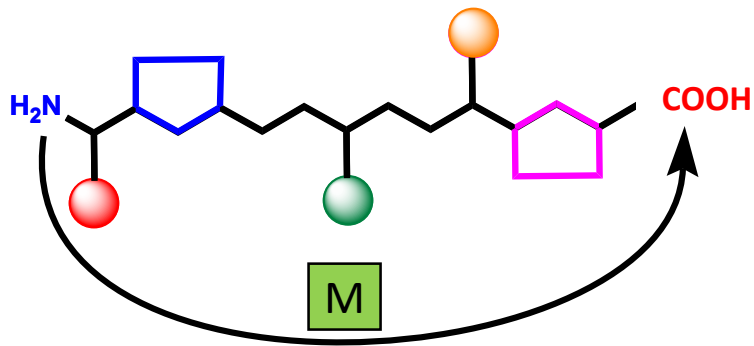
Gyrocycle™ highly modified macrocyclic peptides - effective alternatives to therapeutic peptides

Using Protein Sequence Data

	*	320	*	340	*	360	*	380							
PatG	:	EGEEILVQA	AIKKCQDNNVLIVS	PTGNN	SNESWC	LPAVLP	PGTLAVGAAKVDG	TPCHF	SNWGGNNTKEGILAPGEEILG	: 292					
TruG	:	EGEELLVQA	AVKKCQDNNVLIVS	PTGNL	GECWC	MPAVLP	PGTLGVGAAKVDG	TPCHF	SNWGGNNAEEGILAPGEDVLG	: 345					
TenG	:	VGEEILVKA	AIKKCLDNNILIVAP	VGNNS	NKNWC	LPAVLP	PGILAVGAAKVDG	TPAHF	SNWGGNNTQEGILAPGV	DVLG : 286					
McaG	:	EGEEILVKA	AIKKCIDNNILIVS	PTGNL	GECWC	MPAVLP	PGTLAVGAAKVDG	TPCHF	SNWGGNNGEEGILAPGEDILG	: 347					
LynG	:	KGEELLTQA	AVKKCQDNNILIVS	PTGND	KGECWC	LPAVLP	PGTLAVGAAKVDG	TPCHF	SNWGGNNAEEGILAPGEDILG	: 304					
AcyG	:	VAPDLFARAVKQ	CQDNNMLIVAP	GGNDK	GECWC	IPSIL	PGVITVGAMRDDG	QPFKF	SNYGGEYQNKGVMANGENILG	: 258					
PagG2	:	MAQEIFARAVKQ	CQDSNILIVAP	GGNDK	GECWC	IPSIL	PDVLTVGAMRDDG	QPFKF	SNYGGDYQKQKVMANGENILG	: 259					
OscG	:	IAHDLARAVKNC	QDNNILIVAPT	GNDK	GECWC	IPAIL	PGVLGAGMMKDN	GKPNY	SNWGGNYQHDGILAPGENILG	: 291					
PirG	:	LAHEFIDKAI	RQAQANNILVIAP	GGNDK	GECWC	VPAVLE	NVLTVGAMKDT	GEPFKF	SNFGGKYATQGILAPGENILG	: 246					
TriK	:	IGHEILEKAVRQA	QENNVLIVAPT	GNK	GECWC	LPAILL	PGVMSV	GAMKDN	GQVFKF	SNWGGQYQQQGIAPGENILG : 245					
		e	A64	cqdnN6L66	P	GN1	gecWC6Pa6Lpg	6	vGa	4	G	p	5SN5GG	G66ApGe	6LG

- Many sequences with a variety of origins have the same properties
- If sequences have >30% similarity they are likely to have the same function
- Sequences can be codon-optimised/engineered making them hard to trace back to the source

Using Engineered Enzymes



Heterocyclase



Macrocyclase

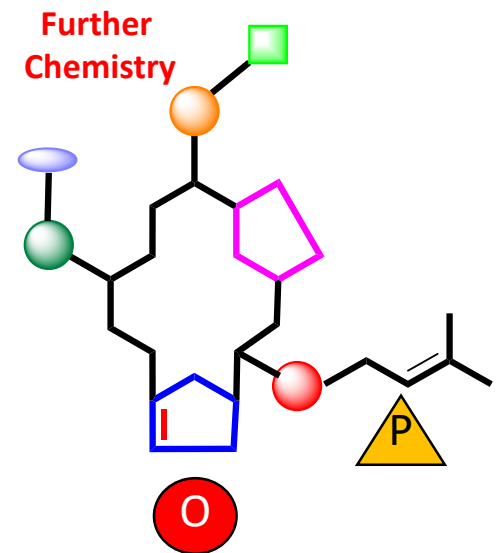


Oxidase



Prenylase

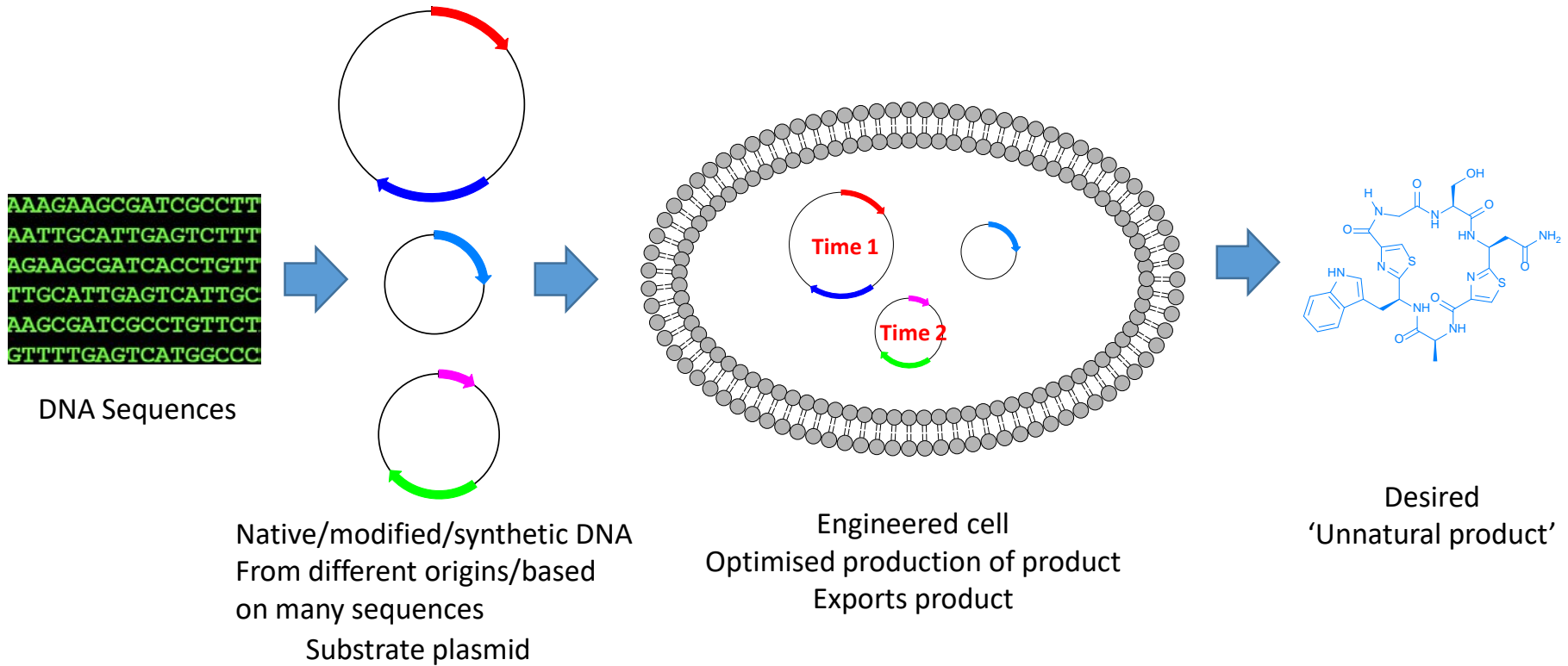
Gyrocycle™ highly modified macrocyclic peptides - effective alternatives to therapeutic peptides



Synthetic Biology/Engineering Biology

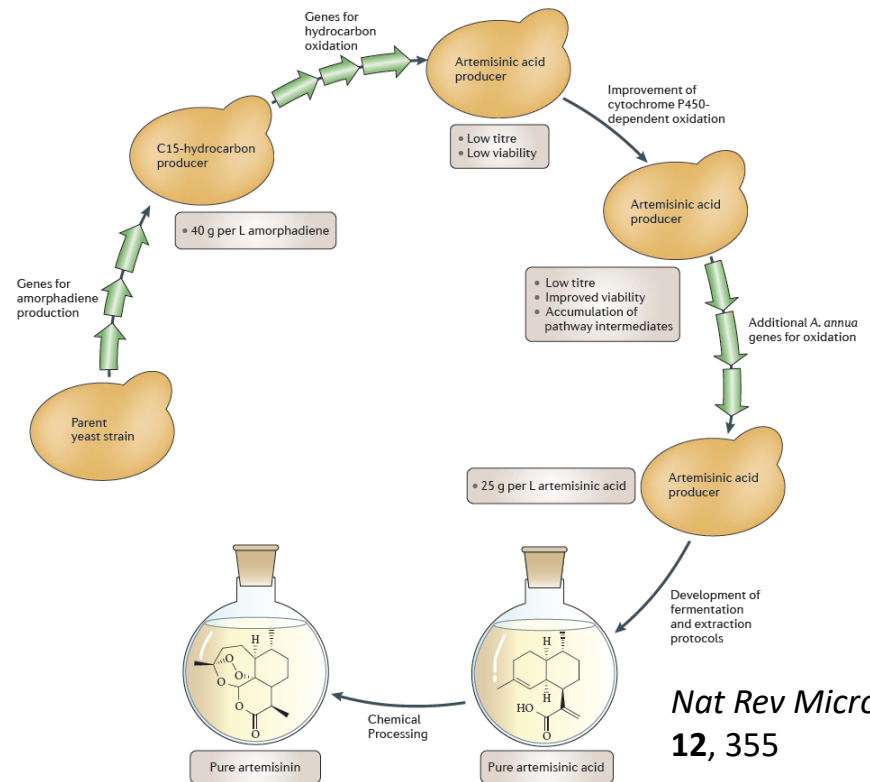
(Design-Build-Test-Analyse-Rebuild)

Together with:



Examples of Application of 'DSI' WHO Essential Medicines Made Using SynBio

- Malaria killed 445,000 people globally in 2016
- Quality-assured artemisinin-based combination therapies are a major part of the strategy against malaria
- SynBio process developed over 7 years with multi million US\$ input
- Produces 25 g/L of precursor that needs synthetic finishing – could meet 25% of global demand



Conclusions

- Many uses of DSI in basic and applied research e.g.
 - Conservation
 - Food safety
 - Vaccines/Pharmaceuticals
 - Industrial enzymes etc.
- Usage of DSI is not straightforward or linear
- Multiple genes sequences used from databases to create final product
 - Some to understand the function of genes
 - Some are incorporated in native or modified form
 - Many genes deposited many years ago
 - **What is the relative contribution of each sequence analysed/used?**