

# ISF OVERVIEW

Vishal Nayak

# WHAT IS ISF?

- ISF = InsertSequenceFeatures
- GUS Plugin to load features located on sequences.
- Designed to deal with
  - Different types of sequence features and feature trees
  - multiple file formats
  - Inconsistency across a single file format:
    - feature names
    - attributes
    - naming conventions
  - many-to-many mappings between features/attributes and GUS tables/columns

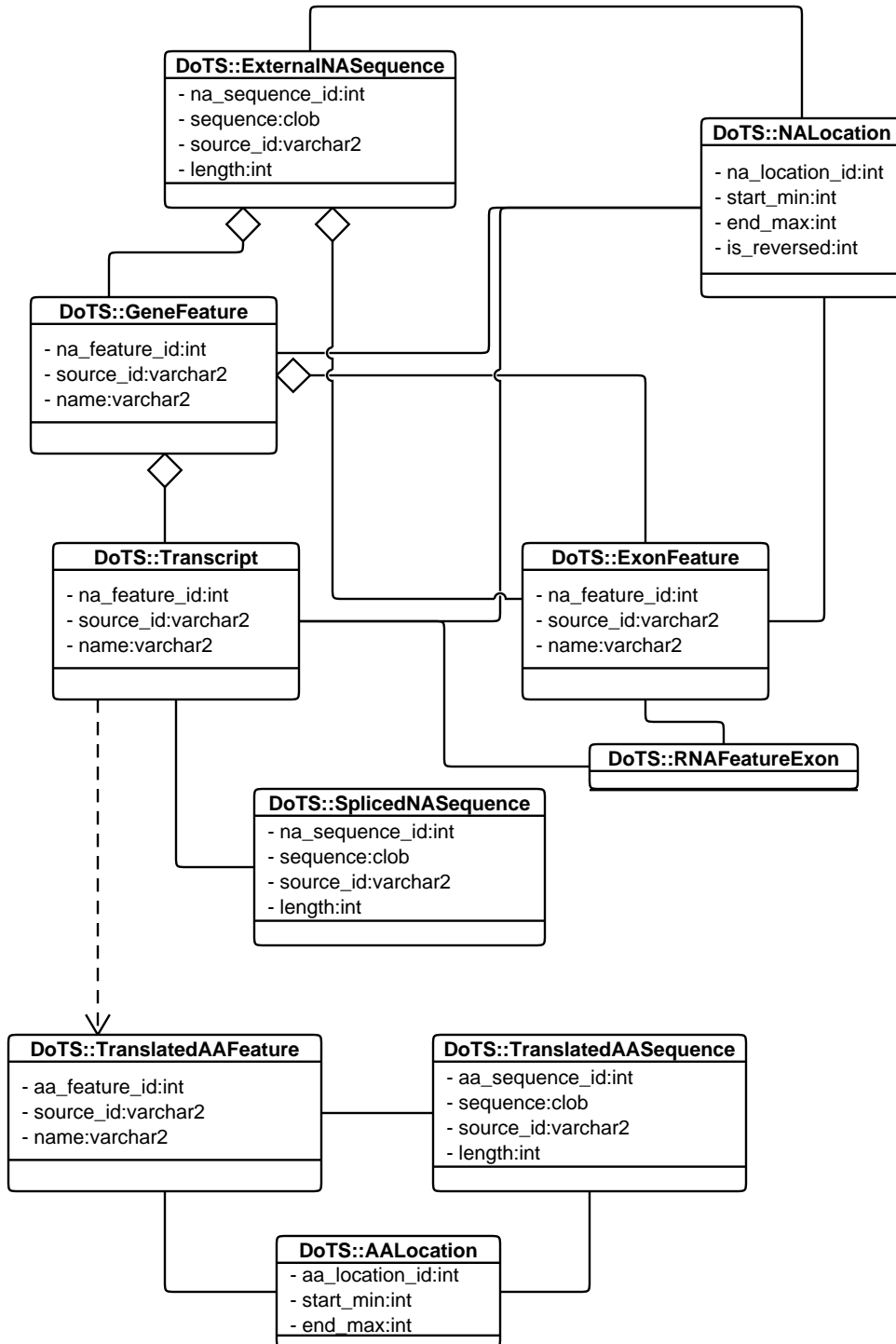


## WHAT IS ISF (CONTINUED)

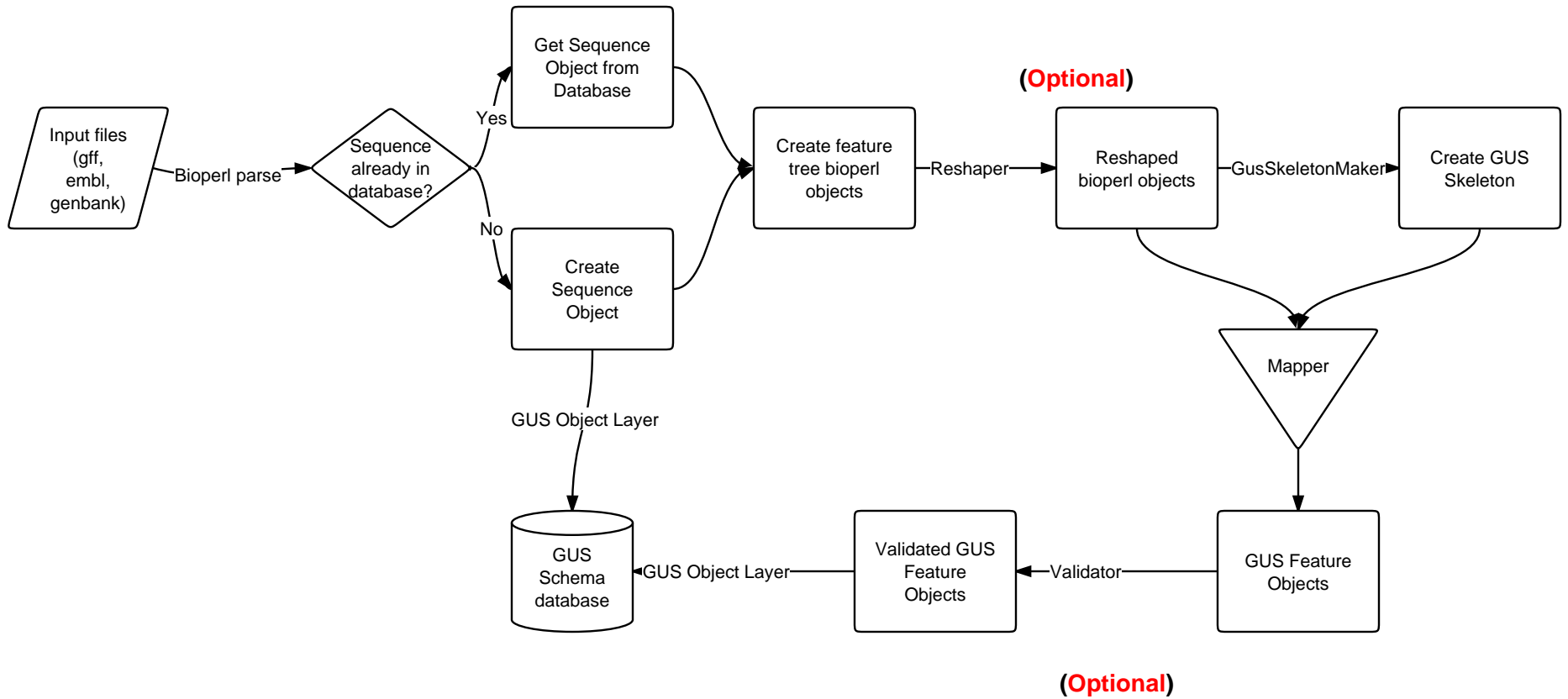
- Supported file formats:
  - genbank, embl, gff2, gff3
- Loads:
  - Features on pre-existing sequences.
  - Or sequences and features simultaneously,
- Examples:
  - gene models, ORFs, isolates, binding site features.



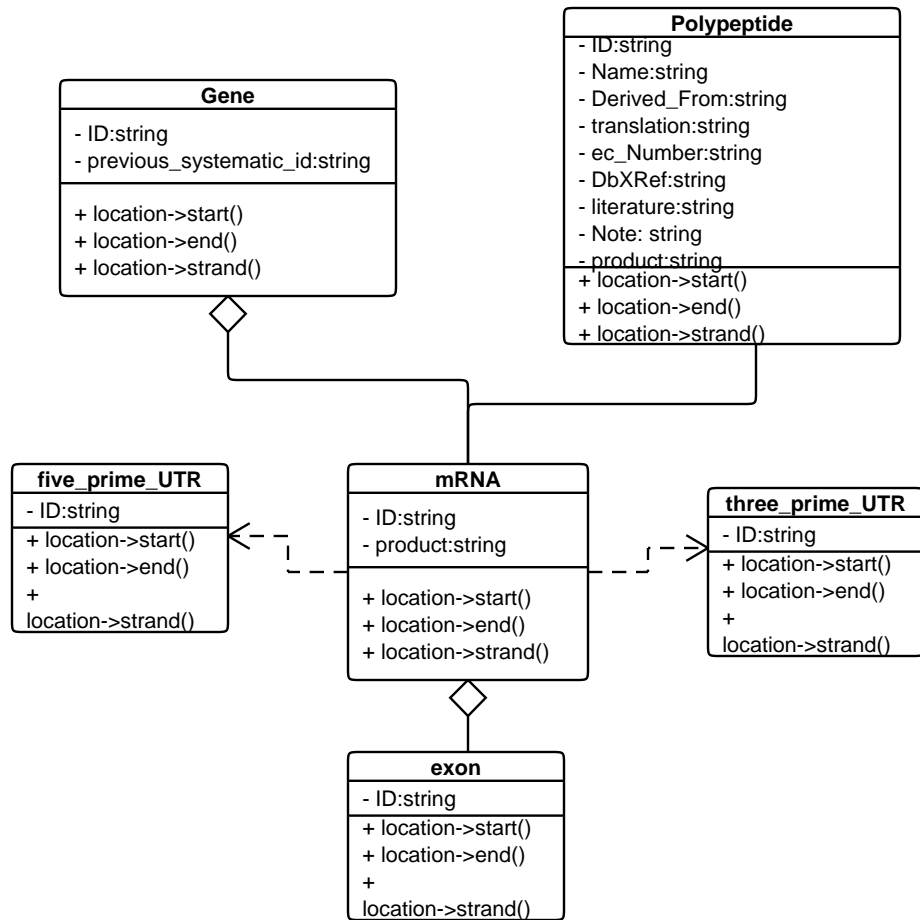
# Gene Models In GUS Schema



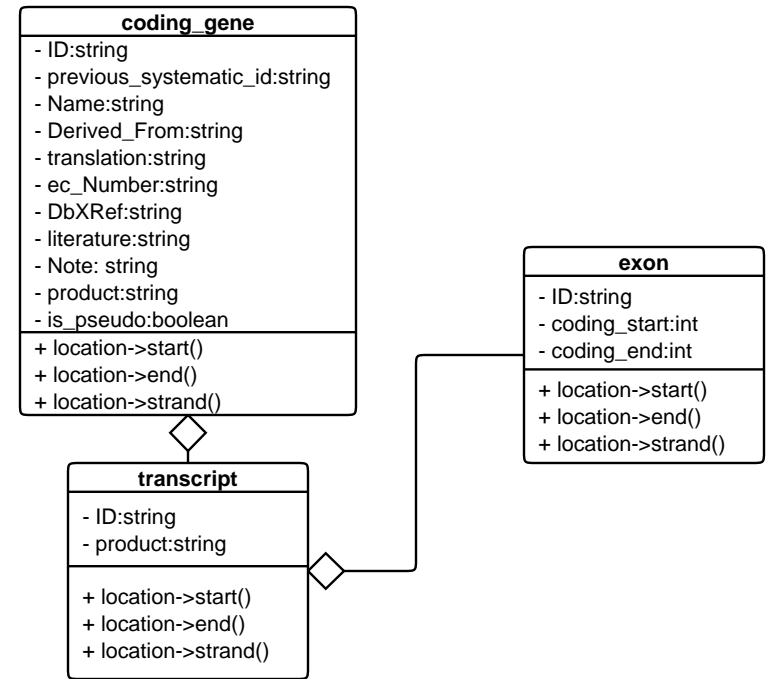
# ISF Conceptual Flowchart



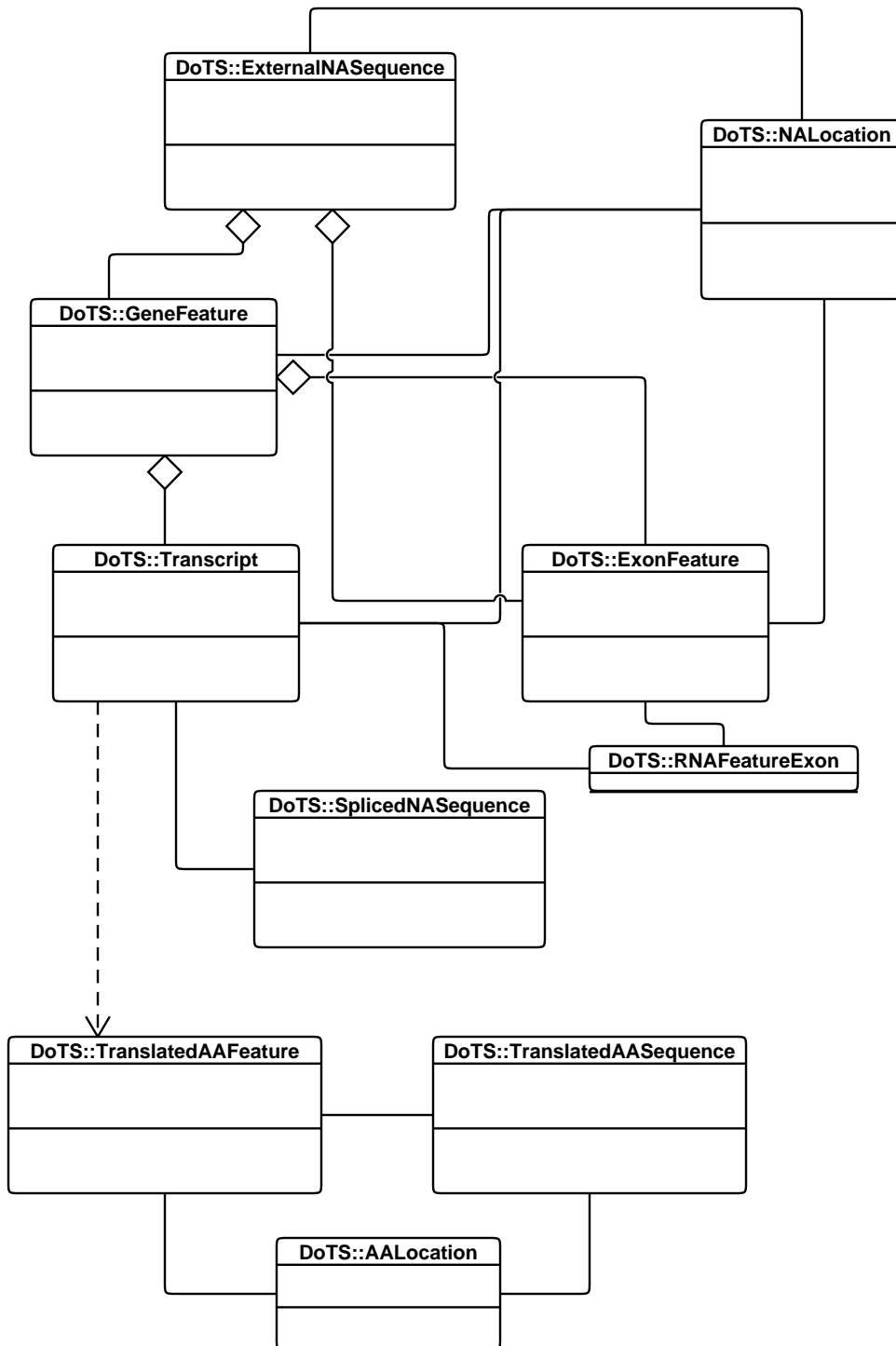
## Sanger GFF Gene Models (BioPerl)



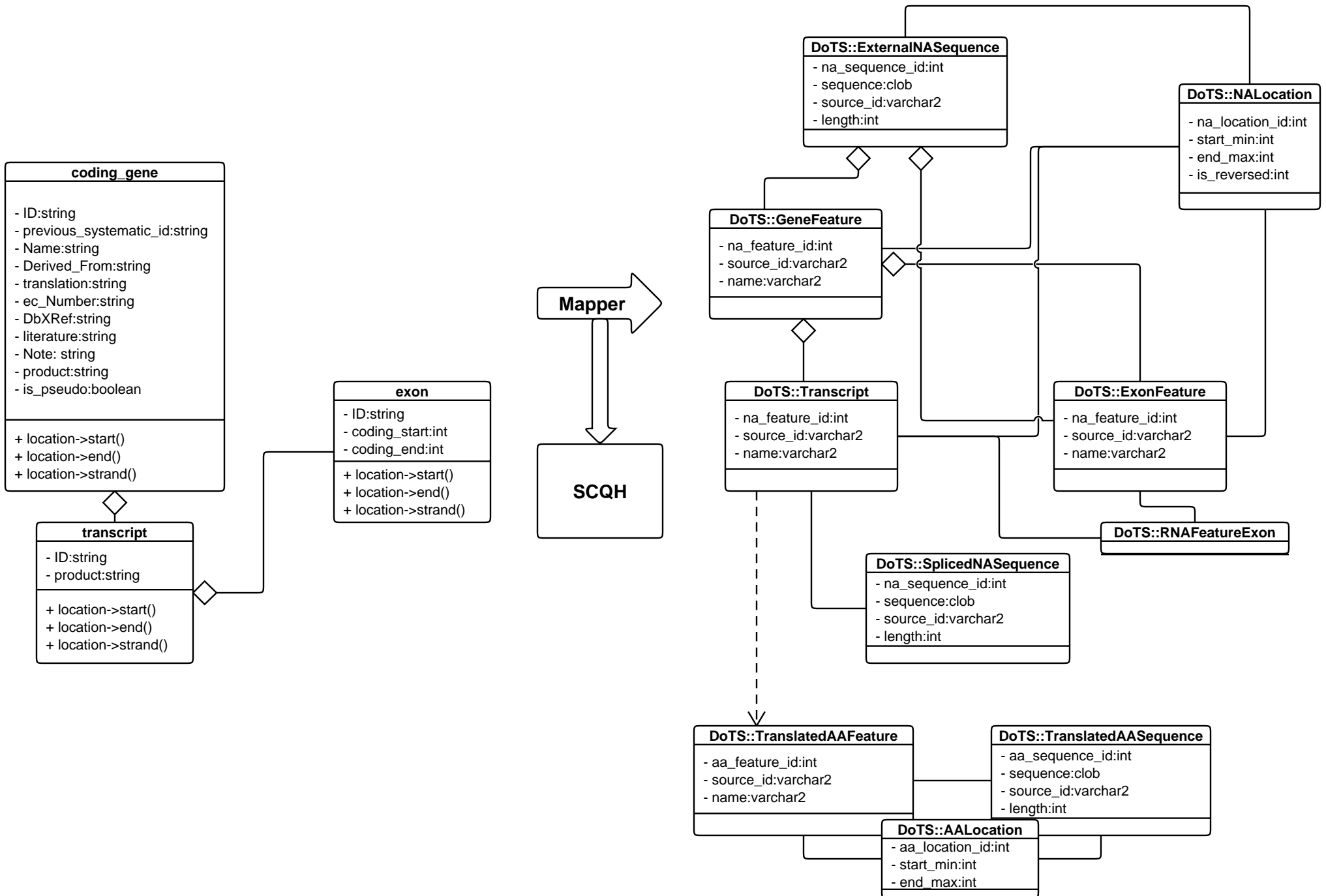
## Reshaped Gene Model (BioPerl)



# GUSSkeletonMaker - Gene Model Skeleton



# ISF Mapper





# MAPPER

- An xml file that is required by ISF
- Maps:
  - Reshaped bioperl objects to GUS objects created by GUSSkeletonMaker.
  - Qualifiers/attributes to corresponding GUS objects/attributes.
- Mapping between reshaped features/qualifiers and GUS objects/attributes might be many-to-many.
- If mapping between attributes is not one-to-one, we need **Special Case Qualifier Handlers**.



# ISF VALIDATION

- Validation is done on feature level (currently, only for coding genes).
- For coding genes:
  - Calculates protein translation and matches it against the provided translation. If translations are not provided, uses the calculated translations.
  - Checks if said gene has a source id.
  - Issues a warning if the sequence contains stop codons in the middle.
- For features requiring validation, a `<validator>` tag has to be specified in the mapper.xml file.



## ISF LIMITATIONS

- Currently unable to deal with gff3 format files directly.
  - Sequences need to be moved to a fasta file and loaded separately.
  - Separate gff3 files need to be merged and preprocessed into a “pseudo gff3” format.
- We are not yet dealing with alternative splicing.
- The validator uses an utility for calculating protein translations that requires all three bases in a codon to identify an amino acid.

