# SEARCH IN MATHEMATICAL DATABASES

Abdou Youssef

The George Washington University

And

The National Institute of Standards and Technology

(DLMF)

# Outline

- Context of the DLMF Math Search Project

- The Project's Short-Term Goals

- Where we are: A Demo

- Technical Issues and Techniques

- Goals and Issues for the Longer Term

# Context of the Math Search Project

- The Digital Library of Mathematical Functions (DLMF) at NIST

    - Web+Book Replacement of the Abramowitz and Stegun Handbook

    - Special functions, Analysis, Functions of Number Theory, Combinatorial Analysis, Numerical Methods, Statistical Methods, …
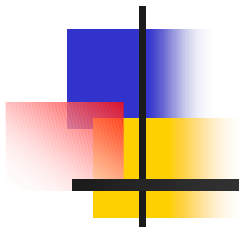
    - DLMF:Mostly Equations – Need Math Search

# Short-Term Goals

n Build a math search system that

1. Understands math symbols & structures

2. Returns equations directly, not just hit-titles

3. Highlights matched equations in documents

4. Understands dialects (Latex, Mathematica, Maple)

5. Provides different search modes (TOC, Index, Free-style search, and Menu-driven search)

# Where we Are

Demo of the Search System

# Sample Queries: understanding math, eq. search & highlighting

| Form | Entry |
|------|-------|
| $\int_0^\infty \sin(\frac{1}{3}t^3 + xt)$ | n int_0^infinity sin((1/3)t^3+xt) <br> n int sin((1/3)t^3+xt) |
| $\sqrt{Ai^2 + Bi^2}$ | sqrt(Ai^2+Bi^2) |
| $\Gamma(\lambda\text{-}\$+\$)$ | Gamma(lambda-$+$) |
| $J_\nu$ or $J_0$ | J_nu or J_0 |
| Ai and J | n Ai and J <br> n Ai and BesselJ |

# Sample Queries: different dialects

- n **BesselJ(nu,z)**
- n **BesselJ(nu, )**
- n **BesselJ( ,zeta)**
- n **JacobiP(n,alpha,beta,x)**
- n **JacobiP( , alpha, , )**
- n **LaguerreL( , , x)**
- n **LegendreP[ , mu , ]**
- n **LegendreP[ , ,sqrt(z)]**

# Search Modes

- One extreme:  (static and limited)
  - <u>Table of contents</u>  (thematic, coarse-grained)
  - <u>Index</u>  (alphabetical, fine-grained)
- Opposite extreme:  (dynamic and unlimited)
  - Free-style search
- Another mode:  (a middle ground)
  - <u>Menu-driven search</u>
    - based on an ontology
    - constrained/standard vocabulary
- Hybrids of the above

# Issues Faced

- Recognizing and Indexing Math Symbols and Structures
- Highlighting Matched Equations (GIF Images) inside HTML Documents
- Development of a Query Language that is Intuitive, Natural, Rich, and Consistent
- Obtaining/Deriving Metadata for Equations
- Development of a Math Taxonomy/Ontology Suitable for Menu-Driven Search

# Techniques: for Handling Math Symbols and Structures

- n  For Recognizing and Indexing Math Symbols and Structures, *TexSN*ize:

    1. Textualization of math symbols (illust.)
    2. Scoping of the various parts of terms/exprs
    3. Normalization of the orders of parts (illust.)

- n  *TexSN*ize the Contents Offline before Indexing

- n  *TexSN*ize each Query before the search

# Techniques: for Equation Search and Highlighting

- n    Create a data model that logically decouples equations from their native documents.

- n    Assign a unique ID to each equation

- n    <u>For Returning Equations directly</u>:
    - n    algorithm that uses a hit list of equation IDs to generate online a document containing the equations

- n    <u>For Highlighting Equations</u>:
    - n    Use the IDs of matched equations to locate the latter in a to-be-displayed document
    - n    add coloring HTML markup to doc before display

# Architecture of the System

- n [Surrogate Files](#)
- n [Indexing System Architecture](#)
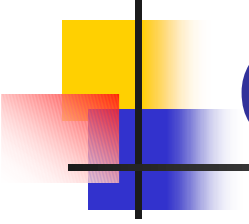- n [Search System Architecture](#)

# Goals for the Longer Term

- Development of a 2$^{nd}$ Generation Math Search System
  - Based on Content MathML+XPath/XQuery
  - More Precise/Expressive Query Language
    - Higher resolution search
    - Keyword search
    - Predicate search
    - Search with term substitution
  - Similarity Search (for Sci. Data Mining)

# Examples of Future Query Types

- n Queries specifying subparts
  - n sin x in a denominator
  - n x-y in a 3rd row of a matrix
  - n $2\pi x$ inside an argument of a function
- n Predicate queries
  - n $z^k$, where $k$ is an integer that ranges from $-4$ to $4$
- n Term-substitution queries
  - n $g(\omega)=z^2 + z + 1$, where $z = e^{i\omega}$
- n Abstraction support and similarity search
  - n x^2 + y^2 = 1 whatever x and y

# Candidate Syntax (Based on 1st Order Logic)

- n /(… sin x …)
- n @(… $2\pi x$ …)
- n z^$k where integer($k) & abs($k)<5
- n x-y in matrix[2,3]
- n x^2 in matrix [$k,$j] where abs($k-$j) < 2
- n $A where matrix($A) &

    (forall $k) (forall $j < $k): $A[$k,$j]>0

- n $S where set($S) & ($\exists$ $x in $S): integer($x) & $x>0
- n x<1 in condition(set)

# Issues that Need to Be Faced

- Canonical Normal Forms of Contents
  - Math equivalences:  ab/c: (a*b)/c or a*(b/c)
  - Notational equivalences: $\int_{a}^{b} \quad or \quad \int_{[\,a\,,\,b\,]}$
  - Distributed definitions
- Uniform Symbolic Notation
- Standard Ontologies
- Development of Metadata
  - Automated extrapolation of metadata
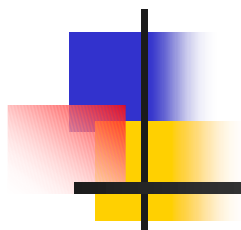  - Manual (by authors and communities)

# Issues to Be Faced (Contd.)

- What Users Need/Want/Prefer
  - What modes of search?
  - What kinds of information?
    - definitions, equations, theorems, proofs, proof techniques, step-by-step evaluations, themes, theories, expositions, etc.?
  - What granularity of retrieval unit?
  - What interactive features?
    - Definition of terms, plotting/computing of matched functions?
- Use of Knowledge of Users' Needs/Preferences
  - More relevant search features & capabilities
  - Better design of the search user interface
  - Better relevance-ranking of search results

# We Are Barely Scratching the Surface

n The Possibilities Are Endless

  n Search + Automated Reasoning

  n Search + Computing + Visualization

# The end