# Protein Identification: Comprehensive Overview

## Database Searching

- Match spectra to sequence databases
- Multiple search engines available
- Statistical scoring algorithms

## Peptide-Spectrum Matching

- Compare experimental vs theoretical
- Fragment ion matching
- Mass accuracy requirements

## Score Calculations

- Probability-based scoring
- Expectation values (E-values)
- Confidence metrics

## FDR Estimation

- False discovery rate control
- Target-decoy approach
- Typically 1-5% FDR threshold

## 1. Database Searching: The Foundation of Protein Identification

Database searching is the primary method for identifying proteins from mass spectrometry data. This approach compares experimental MS/MS spectra against theoretical spectra generated from protein sequence databases. The process involves computational algorithms that systematically search through millions of potential peptide sequences to find the best match for each observed spectrum.

The success of database searching depends on several critical factors: the comprehensiveness of the sequence database, the accuracy of the mass spectrometer, and the sophistication of the search algorithm. Modern proteomics experiments routinely identify thousands of proteins from complex biological samples using this approach.

## Database Searching Workflow

**1** **MS/MS Spectrum Acquisition**

Mass spectrometer generates fragmentation spectra from peptides with specific m/z values and charge states

↓

**2** **Database Selection**

Choose appropriate protein database (e.g., UniProt, NCBI, species-specific databases)

↓

**3** **In Silico Digestion**

Computational digestion of database proteins with specified protease (e.g., trypsin) to generate theoretical peptides

↓

**4** **Candidate Peptide Selection**

Filter peptides based on precursor mass and modifications to create candidate list

↓

**5** **Theoretical Spectrum Generation**

Generate predicted fragment ion spectra for each candidate peptide

↓

**6** **Scoring and Ranking**

Compare experimental spectrum to theoretical spectra and assign scores to rank matches

| Database Type | Characteristics | Best Use Case |
|---|---|---|
| **UniProtKB/Swiss-Prot** | Manually curated, high quality | High-confidence identifications |
| **UniProtKB/TrEMBL** | Automatically annotated, comprehensive | Discovering novel proteins |
| **NCBI RefSeq** | Non-redundant, well-annotated | Model organism studies |
| **Ensembl** | Genome-based predictions | Genomics-proteomics integration |

# 2. Peptide-Spectrum Matching (PSM): The Heart of Identification

Peptide-Spectrum Matching (PSM) is the core process that determines which peptide sequence best explains an observed MS/MS spectrum. This involves comparing the experimental fragmentation pattern against theoretical predictions to find the most likely match. The quality of PSM directly impacts the reliability of protein identifications.

The matching process considers multiple factors including fragment ion types (b-ions, y-ions), mass accuracy, intensity patterns, and the presence or absence of expected peaks. Advanced algorithms also account for neutral losses, isotope patterns, and instrument-specific characteristics to improve matching accuracy.

## PSM Process Visualization

Experimental
**MS/MS Spectrum**
Peak List:
m/z | Intensity

⟷

Matching Algorithm
**Comparison**
Score Calculation

⟷

Theoretical
**Fragment Ions**
Predicted b, y ions

### Fragment Ion Nomenclature

N-terminus   →   Peptide Sequence   →   C-terminus

**b-ions:** N-terminal fragments (retain charge at N-terminus)

**y-ions:** C-terminal fragments (retain charge at C-terminus)

**a-ions:** b-ions minus CO (b - 28 Da)

## Critical Mass Accuracy Requirements:

- **Low-resolution instruments (Ion Traps):** ±0.5-1.0 Da precursor mass tolerance

- **High-resolution instruments (Orbitrap, Q-TOF):** ±5-20 ppm precursor accuracy

- **Fragment ion tolerance:** Typically 0.02-0.05 Da for high-res, 0.5-1.0 Da for low-res

- **Impact:** Tighter tolerances reduce false positives but require excellent calibration

- **Modern standard:** Sub-1 ppm precursor accuracy with <10 ppm fragment accuracy

## Key Matching Criteria:

**1. Mass Matching:** The experimental precursor mass must match the theoretical peptide mass within the specified tolerance. This initial filter dramatically reduces the search space.

**2. Fragment Ion Coverage:** A good match should explain a significant portion of the observed peaks. Typically, identifying 60-80% of major fragment ions indicates a reliable match.

**3. Intensity Correlation:** The relative intensities of matched peaks should show some correlation with theoretical predictions, though this varies by fragmentation method.

**4. Complementary Ion Series:** Presence of complementary b and y ion series strengthens confidence in the identification.

# 3. Score Calculations: Quantifying Match Quality

Scoring algorithms assign numerical values to peptide-spectrum matches, enabling objective ranking and statistical validation. Different search engines employ various scoring strategies, from simple correlation measures to sophisticated probabilistic models. Understanding these scores is crucial for interpreting identification results and setting appropriate confidence thresholds.

Modern scoring approaches consider not just whether peaks match, but also the statistical significance of those matches given the database size, spectrum quality, and other contextual factors. This probabilistic framework allows researchers to control error rates systematically.

## Scoring Methodology

### Cross-Correlation (XCorr)

Used by SEQUEST; measures similarity between experimental and theoretical spectra by sliding them against each other

$$XCorr = \Sigma(Exp \times Theo)_\tau$$

Higher XCorr = better match

### MOWSE Score

Used by Mascot; probability-based scoring considering ion frequency and database size

$$Score = -10 \times \log_{10}(P)$$

Score > threshold = significant

### Hyperscore

Used by X!Tandem; combines matched ion counts with intensity information

### Spectral Probability

Used by MS-GF+; calculates probability of observing spectrum by chance

$$\text{Hyperscore} = N! \times \Sigma(I_i)$$

Factorial emphasizes completeness

$$P(S|peptide) \text{ vs } P(S|random)$$

Direct probability assessment

## Expectation Values (E-values):

- **Definition:** Number of times a match with this score is expected to occur by chance
- **Calculation:** E-value = Database size × P(score by chance)
- **Interpretation:** Lower E-values indicate more significant matches (e.g., E < 0.01 is good)
- **Database dependency:** E-values scale with database size; larger databases → higher E-values
- **Advantage:** Provides intuitive statistical meaning independent of scoring scheme

## Score Interpretation Guidelines:

| Score Type | Good Match Threshold | Considerations |
|---|---|---|
| XCorr (SEQUEST) | +1: >1.9, +2: >2.2, +3: >3.75 | Depends on charge state and peptide length |
| Mascot Ion Score | Typically >30-40 | Identity threshold shown by Mascot |
| E-value | <0.01 (often <0.001) | Lower is better; database-dependent |

| Posterior Error Probability | <0.01 (1% error) | Direct probability of incorrect assignment |

**Important Note:** Raw scores should never be interpreted in isolation. Always consider multiple factors including: spectrum quality, number of matched ions, mass accuracy, score difference between top hits (delta score), and peptide length.

## 4. False Discovery Rate (FDR) Estimation: Controlling Errors

False Discovery Rate (FDR) estimation is the gold standard for controlling identification errors in proteomics. FDR represents the expected proportion of false positives among all reported identifications. Unlike traditional p-values, FDR provides direct control over error rates in the context of multiple hypothesis testing, which is essential when searching thousands or millions of spectra.

The target-decoy strategy has become the dominant approach for FDR estimation. By creating a decoy database of reversed or shuffled sequences, we can estimate the number of false positives without knowing the true identifications. This elegant solution provides empirical, data-driven error estimates.

### Target-Decoy Approach

**1** **Create Decoy Database**
Generate reversed or shuffled sequences from target database (same size, similar composition)

↓

## 2 Concatenated Search

Search against combined target + decoy database simultaneously

↓

## 3 Score Distribution Analysis

Examine distribution of target vs decoy hits across score ranges

↓

## 4 FDR Calculation

At each score threshold: FDR = (2 × Decoy hits) / Target hits

↓

## 5 Apply Threshold

Select score cutoff that achieves desired FDR (typically 1% or 5%)

---

### FDR Calculation Formula

**FDR = (Number of Decoy PSMs / Number of Target PSMs) × 100%**

Alternative formula when using separate searches:
FDR = (2 × Decoy / Target) × 100%

---

**FDR Best Practices:**

- **Multiple Levels:** Calculate FDR at PSM, peptide, and protein levels separately

- **Typical Thresholds:** 1% FDR for high-confidence results; 5% for exploratory studies

- **Decoy Generation:** Use reversed sequences (preferred) or shuffled sequences; avoid scrambled

- **One-Class SVM Alternative:** Machine learning approaches can improve FDR estimation

- **Q-value Reporting:** Report q-values (minimum FDR at which identification is accepted)

- **Avoid Cherry-Picking:** Set FDR threshold before examining results to avoid bias

## Understanding FDR Levels:

### PSM Level

Individual spectrum-peptide matches

**FDR = 1%**

1 in 100 PSMs is incorrect

### Peptide Level

Unique peptide sequences

**FDR = 1%**

1 in 100 peptides is incorrect

### Protein Level

Protein identifications

**FDR = 1%**

1 in 100 proteins is incorrect

**Critical Consideration:** Protein-level FDR is more conservative than PSM-level FDR. A protein identified by multiple peptides has higher confidence than one identified by a single peptide, even at the same FDR threshold.

**Practical Example:** In an experiment with 100,000 PSMs at 1% FDR, approximately 1,000 incorrect identifications are expected. However, if these false PSMs are randomly distributed across proteins and most proteins are identified by multiple peptides, the protein-level FDR will be much lower than 1%.

| FDR Threshold | Interpretation | Typical Application |
|---|---|---|

| 0.1% (0.001) | Very high confidence, few false positives | Biomarker discovery, clinical applications |
|---|---|---|
| 1% (0.01) | Standard high-confidence threshold | Most publication-quality proteomics studies |
| 5% (0.05) | Moderate confidence, more identifications | Exploratory studies, hypothesis generation |
| 10% (0.10) | Lower confidence, many false positives | Preliminary screening, requires validation |

### Integration: Complete Protein Identification Pipeline

Modern protein identification workflows integrate all four components—**database searching** provides candidate matches, **peptide-spectrum matching** evaluates quality, **scoring algorithms** rank results, and **FDR estimation** controls errors—to deliver reliable, statistically validated protein identifications from complex mass spectrometry data.