



Solving the Identity Search & Matching Problem

An overview of the identity search & matching problem
space, and an introduction to Identity Systems' technology

Introduction

A primary requirement of identity subsystems is to query and locate records or information using identity attributes such as Name, Address, and other supporting data.

Anyone who has ever searched a telephone book for a name like Ricky Smith knows that the endless variations and errors in this class of data make for an almost impossible task. (Possible matches include Rick Smythe, Richard Smithe, Smith Richard A., Ricardo A Smith, R. A. Smithe, etc). Computer systems that use names to locate critical information records face the same challenges. The error and variation in names is compounded by the volume of information records being searched against, and the need to perform searches in real-time. In addition, it is typical in such systems that apart from the accidental error and variation – the identity data is subject to *deliberate* abnormal or extreme variation.

On their own, traditional phonetic search techniques like Soundex and NYSIIS, wild-card searches or text search engines, don't address all aspects of the problem space, do not provide matches in ranked order, cannot address words that are in the wrong order (i.e. Lee Kwok Ki and Kwok Ki Lee), and do not handle data from multiple countries/character sets/languages efficiently.

Indexing and search techniques that rely on clean and complete data will struggle with the often incomplete and unpredictable data that must be recorded and searched in such systems.

Reliable, fast identity searches, in addition to accurate matching, duplicate discovery and relationship linking, are fundamental requirements of many computer applications.

The identity search problem is complex – and requires sophisticated tools and techniques to address and solve the problem.

The Identity Search Problem

Names (and other identity data) suffer from unavoidable error and variation.

The variations that occur include spelling, typing and phonetic error; synonyms & nicknames; Anglicization, ethnic, and foreign versions of names; initials, truncation and abbreviation; prefix and suffix variations; compound names; account names; missing words, extra words and word sequence variations, as well as format, character and convention variations.

Variation Examples

Person Names:

Nicknames: William, Bill, Billy, Will

Name Variation: Chris, Kris, Christie, Krissy, Christy, Christine, Tina

Abbreviation/Spelling: Mohammed, Mohd, Mohamad, Mhd, Muhammad

Foreign Versions: Peter, Pete, Pietro, Piere, Pierre

Spelling Variation: Johnson, Johnsen, Johnsson, Johnston, Johnstone, Jonson

Suffix Variation: Smith II, Smith I I, Smith Jr, Smith jnr

Anglicization: De La Grande, Delagrande, D L Grande

Out of Order: Henry Tun Lye Aun; Mr Aun Tun Lye (Henry)

Initials/Order: Frank Lee Adam; A. Frank Lee; Lee Frank

Titles: Dr. Henry Lee, Henry Lee, M.D., Mr. Henry Lee

Company Names:

Town Park Plaza Hotel, Park Plaza, Hotel Plaza, Town Park

John Deer Engg Labs, John Deare Laboratories, Engineering Research Labs c/o John Deere Inc

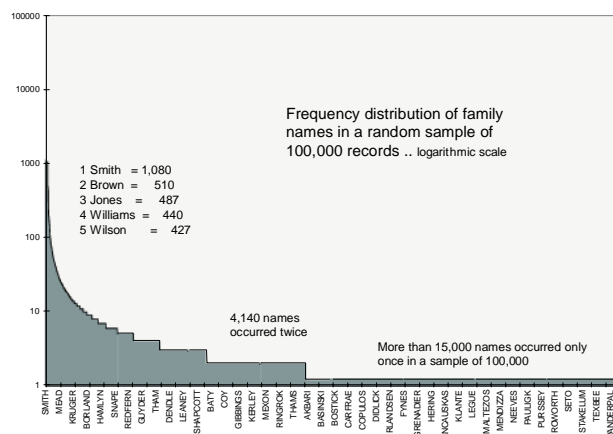
IBM, International Business Machines, Intl. Bus. Mach., I.B.M., IBM Inc.

Addresses:

Jackson Rd. East Hartford
117-2a Jackson Rd East, Hartfrd
2a East Jackson, Hartford
117a Jackson Rd, E. Hartford
Grd. Fl. 192 Aberdeen St Southhead
192/1 Aberdeen Sreet Sth. Hd.
Ground Floor 192 Aberdeen St South Head

Common and Uncommon Words – the Performance Problem

The words we use to label things are chosen from a very different vocabulary than meaningful language. There are no dictionaries, spell checkers or rules for the names of people, places, things or even addresses. The vocabulary in use for people's first names includes in excess of 2,500,000 words in the USA alone, yet as much as 80% of the population may have names from as few as 500 words.



Accurate and high-performance name searching must perform for the uncommon names as well

as for very common words. This is an extremely difficult challenge when a database of 100,000,000 people may contain 100,000 John Smiths, or Juan Rodriguezs or I Main Streets.

International Data

Most large identity databases contain data from multiple languages, countries and cultures which often have different structures, follow different parsing rules, and have different variation characteristics. Also, if transliteration, Romanization, character set conversions and other such transformations are employed, a new class of error and variation is introduced.

Names are Truly Many-to-Many

It is obvious that two people or companies, or products, can have exactly the same name. It is also obvious that, even ignoring error and variation, people places and things have more than one name:

- People have maiden names and married names;
- People have aliases and professional names;
- Companies have registered names, trading names and division names;

- Places have several addresses, on two separate streets, old addresses, billing addresses, postal addresses etc.
- People and places can have names in more than one language.
- The relationship between a name and that which it names is quite naturally a true "many-to-many relationship."

Indexing these "many to many" relations requires careful design in the majority of today's search applications. Search techniques that depend on two fields, one for "name" and one for "maiden name", or "alias" are not good. When we are searching for a person's name, or address we do not know which "role" it plays. We do not know if it is a birth name, married name or maiden name, we do not know if it is a current or prior address. In order to address this problem effectively, it is necessary to have several keys or index entries pointing to the same identity.

The Standardization, Cleansing, and Parsing Dilemma

There are multiple Marketing objectives that have conflicting needs in the way that name and address data is captured, stored, and used. In many marketing systems, this conflict has not been recognized, leading to a bias in one area, and a less than satisfactory solution in another.

Examples of such conflicts are:

The **address** most useful for reaching the prospect and fostering a good relationship is the one the **prospect** provides; the address most useful for achieving a cheap mailing rate is the one the **post office** provides.

The **name** most useful for matching is an unformatted **original**; the name most useful for personalization will be **formatted and enhanced**:

The **address** most useful for matching or creating household groupings is the unformatted or partially formatted **original**; the address most useful for statistical analysis will be **formatted and enhanced**.

Therefore, to achieve high-quality searching and matching of identity data, and to support good prospect/customer relationships, it is necessary to store the original, unenhanced data. In addition, to support marketing needs such as personalization or postal enhancement, it is necessary to store the formatted, enhanced data (or a derivative, e.g., postal barcode), or view these needs as 'output' processes to be applied to the data as it is taken out of the system.

Furthermore, when focusing on searching and matching addresses, it is important to understand exactly how the address is to be used. Traditional address matching methods rely on the cleaning and formatting of an address to build a stable "match-code." However, many address types, such as the location for a particular utility service or a record for land titles, cannot easily be "formatted."

Another aspect to consider in relation to mailing addresses is the source of the address. Increasing commercial globalization and Internet commerce has led to a much greater diversity of addresses reaching our systems. In current times, for example, it is not uncommon for a bookstore in the USA to deliver to a client in Japan.

The format of mailing addresses in some countries is relatively simple and well understood. However, many other countries have much more complex addresses and less rigorous standards for their writing. In addition, localized web sites often request information formatted in the home country's address format rules, which further complicates the issue and results in even messier data.

When designing name and address systems where the system has to operate in a variety of countries, or upon databases that are partitioned for different countries, it is necessary to do all data acquisition and database design in such a fashion as to minimize the differences between countries. This is the model on which many web-based systems must be based. Any organization dealing with addresses from multiple countries cannot rely on searching and matching processes that require that the addresses be cleaned and formatted.

Popular but Less Successful Techniques

Cleaning, Standardizing or Enhancing Names and Addresses

While the concept of cleaning, enhancing or otherwise standardizing names and addresses to make search and matching easier is very enticing, there is a logical conflict between the idea of “cleaning” and the “unpredictability” of identity data. Incorrect assumptions made by formatting and cleaning routines can corrupt good data, and records rejected by the cleaning process can mean lost opportunities to match.

Cleaning and enhancement of identity data may be a valid process for preparing names and addresses for some marketing functions, however such pre-processing does not solve all of the data matching problems, particularly when the data is unreliable or the need to find a match if one exists is critical.

Text Searching

The idea that text retrieval packages can successfully be used for name search applications is only believed by users who are unaware of the “good” data that they miss. Even when “Full Text” indexes have phonetic algorithms, or “expert” rule bases for Name searches the indexing mechanism will result in an inefficient process. Imagine the cost of finding all the index values for the records containing John and then joining them with those that contain Smith to discover the subset John Smith. In actuality, many text search packages fail to recognize that Bill Smith and William Smith could be the same person, let alone some obscure abbreviation such as Wlm Smith.

Typical free text indexing techniques do not allow high performance or high quality retrieval of data containing names. Text searches are not for the serious searcher dealing with identity data.

In those applications where there is a need to index, search and classify both identity and non-identity data in free text or unstructured data – a combined solution is necessary. The non-identity data can be handled by text retrieval and categorization engines – but the identity data must be handled by a sophisticated identity search & match product that recognizes, finds or discovers the “name phrases” in the text and indexes them with the same specialized techniques that are necessary for quality and performance in any other name search system.

Wild-Card Searches

Wild-card searches do overcome SOME error and variation in the name. That is why users believe in them. If the users knew what data they were missing they would not believe in them. Programmers believe in them because they learned to program them in their computer science classes at university. In reality they only work if the searcher correctly guesses the right characters to include or exclude, and the database had no error in the characters actually used. They cannot find all the relevant records, do not address nicknames and abbreviations, nor the fact that different records have different errors.

Wild-card searches normally return too many irrelevant candidates and will always miss many of the relevant candidates. Wild-card searches are not for the serious searcher.

Match-Codes

A Match-code is an entity key built from a combination of the entity's attributes. The Match-code for: John Smith dob: 22 Feb 1979 2/234 33RD St., New York, NY 12345, might be: SMITHJ79NY. Match-codes require that the entity is first strictly formatted into its pieces; that all pieces used are in the “stable” order; and that there are no errors in the pieces used.

Once each source record has a Match-code, “like” Match-codes can be used to bring together similar records for more intense matching. Some match-coding aims to build a unique key for entities such that search and match can be done in one step.

Match-codes can also be used in real-time searching. Depending on their granularity, they can make getting to a matching record very quick. The search width can also be expanded by truncating attributes (e.g. search for Family Name + Initial + Birth Year) without including the State Code).

The problem with Match-codes is that to be really effective they require error-free data, which is why they depend so heavily on cleaning and formatting routines to prepare the data. A missed match could be caused by something as simple as a transposition during data entry (213 Main Street vs. 123 Main Street).

Typically Match-Codes do find “correct records,” but they frequently miss other good candidates. If there is no ability to overcome error and variation in the name parts, all records with such error and variation will be missed. Missing word, extra word and word order variations are often missed unless the searcher permutes the search criteria or the batch process is run in multiple passes.

Match-Codes fail if any one piece of the data used to build the code is not identical. Match-Codes are not for the serious search or matching application and definitely not for the critical search applications or large volume systems.

N-Gram Indexing

An n-gram is a set of “n” consecutive characters extracted from a word or code. Typical values for “n” are 2 or 3. These extracted n-grams are subsequently indexed for all names or addresses in the database. At search time, the idea is that words or codes that are similar between the search and file data will have a high proportion of n-grams in common. N-grams are particularly

well suited to string and text searching; however, unless supported by extensive rule bases for phonetic and synonym variation, as well as for noise words, they do not readily overcome the typical error and variation found in identity data, nor do they easily scale to very large data volumes.

Yesterday's Soundex's, NYSIIS, and other simple Phonetic Algorithms

In the early 1900's the Russell Soundex, was developed to provide a stable manual filing code for the USA Census documents. The development of this algorithm for encoding the last name of a person was based upon phonetics and certain classes of typical spelling and filing errors. It was a simple set of rules to convert a last name word into a four, five or six digit number that had a high probability of being the same for two words that were variations of each other.

Since then, many Algorithms with a similar objective have been developed and modified. In 1969 the New York State Identification Intelligence System project evaluated the then popular algorithms. This included evaluation of last names on New York State criminal records to see how effective algorithms were at overcoming the error and variation. To this they used records that were previously paired using fingerprints. They evaluated such algorithms as: Soundex and many of its variants; LA County Sheriff; Consonant coding; Phonic standard and extended; Michigan Lien; and several Extract list based systems. The end result of their project was to define a popular algorithm known as NYSIIS which they proved was better optimized for their data at that point in time.

At Identity Systems, we have conducted many such exercises on data from many customers and countries and have found that, even in one country, **no single algorithm is optimum for all naming data.**

Fundamentally the choice of the optimum word stabilization algorithm depends upon the source

of the data and the frequency distribution of certain classes of error in the data. For example carefully written and captured Criminal Records have a very different error distribution than urgent verbal requests for information over a radio. While such stabilization algorithms can be a critical piece of a name search engine, the algorithms of the past are not suitable for use as database keys with today's data or volumes. Typically they either cause too many records to be found, or they miss too many relevant records.

Where stabilization algorithms are used, the choice of the right algorithm can be critical and must be based upon the true distribution of the error and variation in both the population of file

data, and the population of the search data. The objective of word stabilization algorithms should be limited to retain as much of the original words form as is commensurate with the objective of "not missing valuable records."

However, the problem of phonetics and spelling is only one part of the general name search problem – and as such a complete name search solution cannot be solely dependent on just a phonetic algorithm. **The design of suitable database keys to maximize quality is a separate and much more complex exercise.**

Identity Systems' Approach

The ideal identity search solution:

- Overcomes the error and variation in the data
- Maintains performance even in extreme volume situations
- Finds all valid candidates without generating false matches.

This implies:

Intelligent and scalable **algorithms**, which, through the use of rich keys and search strategies, return all of the candidates an expert user would consider as being the same as the search criteria.

These algorithms must be able to cope with data from the real world. This includes data that is not formatted or cleaned or which contains missing, extra, truncated, out of order, non standard, or noise characters/words, initials, abbreviations, nicknames, numbers, codes, and concatenations.

The algorithms need a customizable rule base to incorporate the knowledge of the expert

user, and a default population rule base in the case where the user is not that experienced.

The algorithms require phonetic and orthographic correction functionality, to address spelling and typing errors.

Intelligent matching routines must be available and tuned to mimic the expert user making a choice as to which candidates are the correct matches. Such matching routines need to take into account all of the error and variation in the identities' attributes, as well as weighting the attributes as the user would.

The Algorithms must work well regardless of the country of origin and language of the data and must insulate the application developer from the differences between country and language.

Identity Systems software is designed and built on these premises and is being used by over 500 commercial and governmental organizations around the world.

A Representative List of Identity Systems Customers

USA Users include:

Bank of New York
Internal Revenue Service (IRS)
Lockheed Martin
General Motors
CitiCorp
VISA
Equifax
Experian
General Reinsurance
Insurance Services Office
The Hartford
Cendant (Resort Condominiums International)
L. L. Bean
State of California (several systems)
LA County Consolidated Criminal History System
New York State Office of Child and Family Services
New York City Department of Finance
Omeda Communications
AT&T
Sprint
FedEx
Pacific Gas & Electric
Purdue University

International Users include:

Canada Customs & Revenue Agency
Citizenship & Immigration Canada
Australian Bureau of Criminal Intelligence (ABCI)
Brønnøysund (Norway; Government))
Bausparkasse Schwabisch (Germany; Bank)
Benkar Tuketici Finansmani (Turkey; Credit Info)
Serasa (Brazil; Credit Info)
Institute of Medical & Veterinary Science (Australia)
Dun & Bradstreet (D&B) Global
Bradesco Seguros (Brazil; Insurance)
ING Canada
Zurich Financial Services
European Patent Office
Australian Taxation Office
Hong Kong Customs
UK Home Office - Offenders Index
UK Customs & Excise
Czech Security Information Services
Readers Digest (US, Mexico, Brazil, Czech Rep.)
Australian Gas & Light
Electricite de France
British Telecom
Taylor Nelson Sofres (UK, Market Research)

A Representative List of Identity Systems OEM Partners

Siebel (Americas, Europe, China, Japan, Korea, Asia/Pac)
CSC (North America, Europe)
Orion Systems (Asia/Pac & North America)

Ericsson (Europe)
Lightbridge (Global)
Kana (Global)

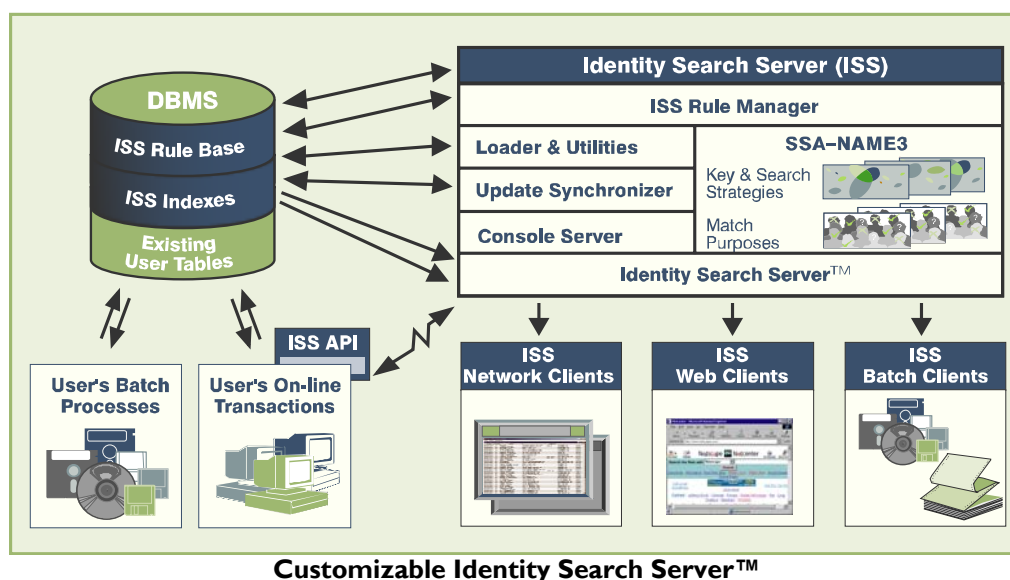
Identity Systems Products

There are three products:

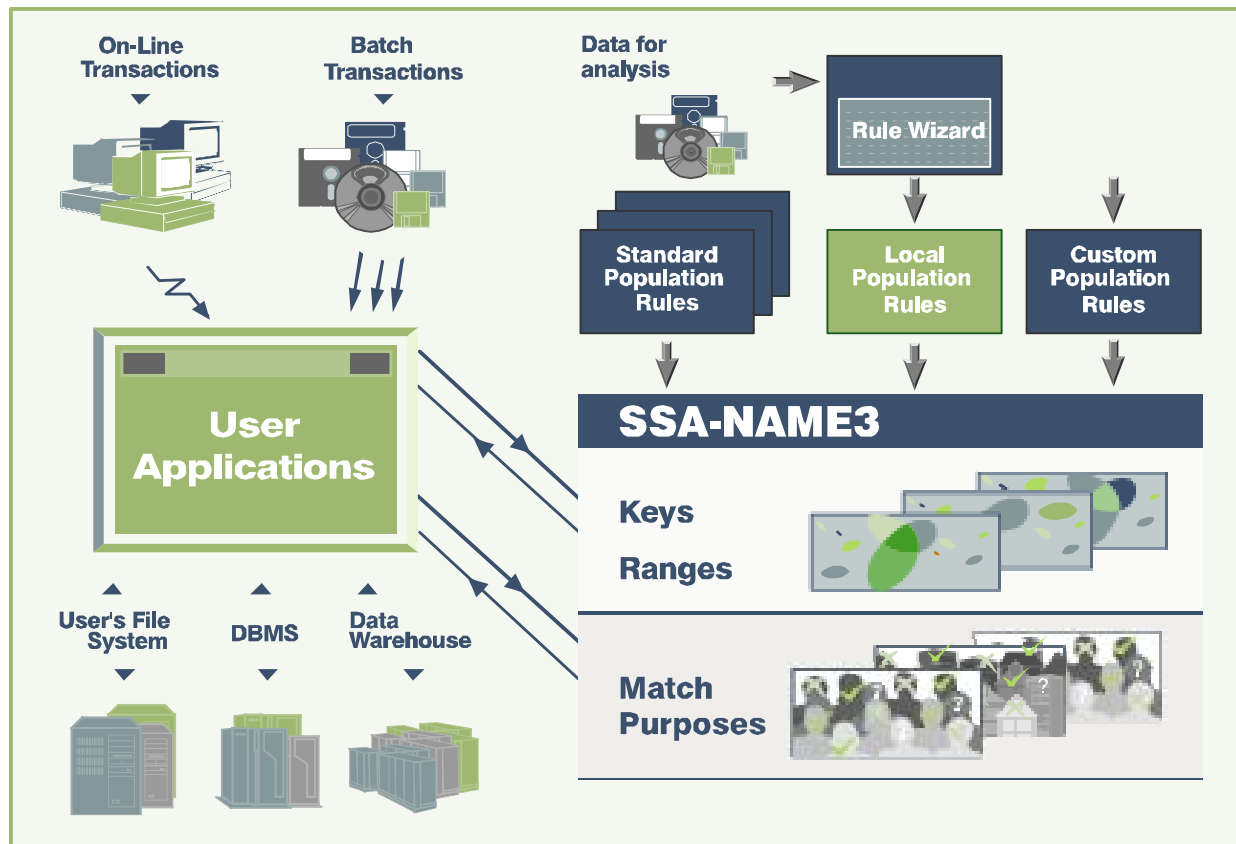
- **Identity Search Server (ISS):** Online / Batch
- **SSA-NAME3:** Core SDK/API (included in other products)
- **Data Clustering Engine (DCE):** Batch

All these products can search, match, and rank identity **data of any quality or format** and do not require any "cleaning" or "scrubbing" of the source data nor is the source modified in any way by Identity Systems products. **All countries/languages** are supported, both singularly and in combination. All products can be used to centrally index and simultaneously search, match, and group identity data from disparate sources such as from Current Sales, Legacy, and External databases. An overview of each product follows.

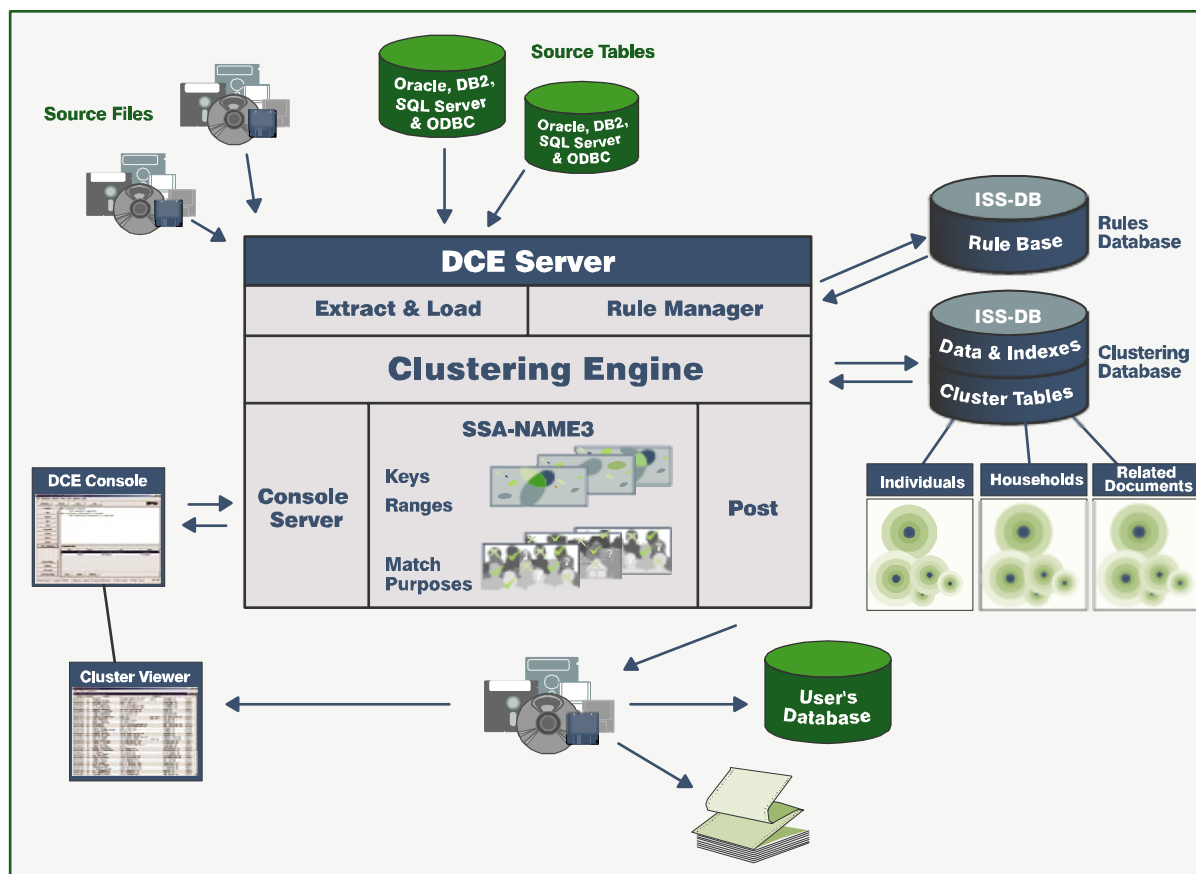
The Identity Search Server (ISS) is the product that provides **on-line and batch searching, matching and duplicate discovery** for all types of identification data stored in relational databases (**DB2/UDB, Oracle, and SQL/Server**) on OS/390, NT, AIX, HP/UX, Linux, & Solaris platforms. High-performance indexes are automatically maintained **without changes to existing application programs**. Custom search clients can be developed using an easy-to-use API. ISS uses an **n-tier** architecture, making it the ideal solution for organizations implementing the latest technologies, including web-based applications. With ISS, the capability to **centrally** index identity information from **disparate** source tables, databases, and computers is particularly advantageous since the central index can then be used to search, match, rank, group, and maintain all the data **simultaneously** and in **real time**.



SSA-NAME3 is Identity Systems **core technology** and included in all the other Identity Systems products. It is a system development kit that enables organizations to build business application programs to **search, match, rank, and analyze** records about people, companies, products, addresses and various other identity data. It can be executed on most any platform and operate on data stored in any database. The software contains algorithms for computing 'fuzzy' search keys and strategies, as well as algorithms for complex matching of identity data.



The Data Clustering Engine (DCE) is a **stand-alone, batch data grouping and investigation engine** for all forms of identification data and can execute on AIX, HP/UX, Linux, Windows NT, Solaris, & Compaq Tru-64 Unix platforms. The DCE performs a thorough analysis of the relationships between people, companies, addresses, products and other entities. This powerful software is used by numerous leading corporations and institutions for **de-duplication, file matching, householding, fraud investigation**, and various other analyses. It **does not require programming** and is highly scalable for the processing of large and extreme data volumes.



Identity Systems

Identity Systems is the pioneer in enabling organizations to build and maintain high-quality identity data search and matching software solutions. Identity Systems have been in this area of specialization for more than 20 years. Identity Systems now has over 600 clients worldwide who rely on its robust enterprise-wide software. The Company has offices in the United States, the United Kingdom, Australia and Brazil. Identity Systems became a wholly owned subsidiary of Nokia in 2006.

Company History

Identity Systems was originally established in Australia in 1987 as Search Software America (SSA) with the purpose of formally developing and marketing its founders experience in building identity search and matching systems. This experience had shown that, while the significance of identity matching in different systems varies considerably, there has always been the same basic set of concerns:

- **Performance problems arise because the most frequently occurring names are also usually those upon which searches are most often performed.**
- **Traditional phonetic algorithms create poor response time and prevent the user from locating a match among many candidates. Traditional match-codes fail to find matches when data is invalid or incorrectly parsed.**
- **True phonetics is only a subset of the errors in names, addresses and other identity data.**

Some of our projects emphasized the need to quickly achieve a match, if there was one. Others placed their emphasis on proving that there was no match at all. One project presented the unusual opportunity to empirically develop and modify an algorithm designed to solve phonetic, orthographic, and Anglicization problems in more than 2,000,000 hand-written credit records. During the project development activity, some 300,000 computerized matches were compared to manual matches done by expert searchers. Whenever the searcher found data not found by the system, the algorithm was refined.

Another project involved the re-processing of 25 million records of international data from virtually every country in the world, where it was known that at least 99% of the records were in fact pairs about the same person. The project also demanded a performance breakthrough since the design objective was to develop an identity search system to handle over 30,000 inserts per day and to support a 50,000,000 record on-line database. Based on the fact that the project's purpose was to identify the records that were not pairs and the population was smaller than the error rate in the data, the successful solution required considerable research.

The experience gained from such projects and Identity Systems work with customers around the world for the last 20 years, has been incorporated into ongoing development and maintenance of Identity Systems product suite. Although the identity search and matching problem may never be completely solved, Identity Systems enables organizations to reach unparalleled levels of quality and performance.

Contact Us

Visit us on the web at www.identitysystems.com or at the contacts below:

USA, North and South America

USA Sales

USASales@identitysystems.com

tel. +(203) 698 2399

Identity Systems

1445 East Putnam Avenue

Old Greenwich CT 06870, USA

UK and Europe

UK Sales

UKSales@identitysystems.com

tel. +44 (118) 944 9688

Identity Systems

9 Headley Road

Woodley READING RG5-4JB, UK

Australia and Asia

Australia Sales

AUSSales@identitysystems.com

tel. +61 (02) 9571 1300

Identity Systems

Nokia House

Level 6, 19 Harris St

Pymont NSW 2009, Australia

Argentina

Levi Minond

eguiterman@leviminond.com.ar

tel. +54(11) 343 2005

Levi Minond

Estudio LmyA, Piedras 575 piso 1 "D"

(1070) Buenos Aires, Argentina

Brazil

Brazil Sales

USASales@identitysystems.com

tel. +0800 891 7052

France

France Sales

FRSales@identitysystems.com

tel. +33 (08)10 85 93 62

Germany

Identity Systems

desales@identitysystems.com

tel. +49 69 677 33 462

Identity Systems

Lyoner Straße 26

D-60528 Frankfurt/Main, Germany

India

Identlogic Systems

sanjayb@identlogic.com

tel. +91 (22) 2496 0033

Identlogic Systems

Hansraj Pragji Building

83C Moses Road

Worli, Mumbai 400018, India

Portugal

Datinfor

ContactCenter@datinfor.com

Datinfor

tel. +351 219 426 600

Quinta do Património, 5

2685-005 Sacavém, Lisboa, Portugal

South Africa

Bateleur Software Solutions

sam@bateleur.co.za

tel. +27 (11) 691 1600

Bateleur Software Solutions

Homestead Park, Block A

65 Homestead Ave

Bryanston, Sandton, Gauteng 2120, South Africa

Turkey

Sales

info@komtas.com

tel. +90(212) 288 47 30

Komtas

Yildiz Posta Caddesi

Dedeman Ishani Kat 10 No.48/23

Esentepe 34349 Istabul