# An Adaptive Hybrid Pattern-Matching Algorithm on Indeterminate Strings

W. F. Smyth[1], Shu Wang[1] and Mao Yu[1]

[1]Algorithms Research Group
Department of Computing & Software
McMaster University, Canada
email: smyth,shuw@mcmaster.ca

The Prague Stringology Conference 2008

# Outline

# Outline

# Outline

# Outline

# Outline

# Outline

**Introduction**
Fundamental Algorithms
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

# Outline

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Regular Pattern Matching Algorithms

Over the last several decades, dozens of regular
pattern-matching algorithms have been proposed.

- Window shifting: KMP [KMP77], BM [BM77], FJS [FJS06], etc.
- Bit-parallel: Shift-Or [Döm68, WM92, BYG92], BNDM [NR98], etc.

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Regular Pattern Matching Algorithms

Over the last several decades, dozens of regular pattern-matching algorithms have been proposed.

- Window shifting: KMP [KMP77], BM [BM77], FJS [FJS06], etc.
- Bit-parallel: Shift-Or [Döm68, WM92, BYG92], BNDM [NR98], etc.

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Indeterminate Pattern-Matching Algorithms

Intuitive approach: Modify existing regular pattern-matching algorithms to do indeterminate pattern-matching.

- Shift-Or: Can be modified to indeterminate pattern-matching easily, with the same speed of regular pattern-matching.

- iBMS: A very fast indeterminate pattern-matching algorithm based on BMS has been proposed in [HSW06b].

- iFJS: An indeterminate pattern-matching algorithm based on modified FJS (cut-off border array) has been proposed in [HSW06a].

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Indeterminate Pattern-Matching Algorithms

Intuitive approach: Modify existing regular pattern-matching algorithms to do indeterminate pattern-matching.

- Shift-Or: Can be modified to indeterminate pattern-matching easily, with the same speed of regular pattern-matching.
- iBMS: A very fast indeterminate pattern-matching algorithm based on BMS has been proposed in [HSW06b].
- iFJS: An indeterminate pattern-matching algorithm based on modified FJS (cut-off border array) has been proposed in [HSW06a].

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Indeterminate Pattern-Matching Algorithms

Intuitive approach: Modify existing regular pattern-matching algorithms to do indeterminate pattern-matching.

- Shift-Or: Can be modified to indeterminate pattern-matching easily, with the same speed of regular pattern-matching.
- iBMS: A very fast indeterminate pattern-matching algorithm based on BMS has been proposed in [HSW06b].
- iFJS: An indeterminate pattern-matching algorithm based on modified FJS (cut-off border array) has been proposed in [HSW06a].

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

## Indeterminate String

A (regular) string $x$ on $\Sigma$ is a finite sequence of letters drawn from $\Sigma$. Two letters $\lambda, \mu \in \Sigma$ are said to *match* ($\lambda \approx \mu$) iff $\lambda = \mu$. Consider any specified subset $S = \{\lambda_1, \lambda_2, \ldots, \lambda_j\}$ of $\Sigma$, $j \geq 2$. We introduce the idea of an indeterminate letter $\lambda = \lambda_S$ with the property that it matches every element of $S$ (but no other letter); we write

$$\lambda \approx \lambda_1, \ \lambda \approx \lambda_2, \ldots, \lambda \approx \lambda_j.$$

Given two subsets $S, T$ of $\Sigma$, $|S| \geq 2$, $|T| \geq 2$, and indeterminate letters $\lambda, \mu$ associated with $S, T$ respectively, $\lambda \approx \mu \Leftrightarrow S \cap T \neq \emptyset$. Given two indeterminate strings $x$ and $y$, $x \approx y \Leftrightarrow (|x| = |y|) \wedge (\forall i \in 1..|x|, x[i] \approx y[i])$.

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# Outline

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

**The Knuth-Morris-Pratt Algorithm**
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

## The Knuth-Morris-Pratt Algorithm

- A well-known linear time pattern-matching algorithm.
- Based on border array calculation.
- However, not very fast in practice.

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

**The Knuth-Morris-Pratt Algorithm**
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The KMP Algorithm - 1



match

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

**The Knuth-Morris-Pratt Algorithm**
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The KMP Algorithm - 2

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

**The Knuth-Morris-Pratt Algorithm**
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The KMP Algorithm - 3

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

## The Sunday Adaption of Boyer-Moore Algorithm

- A simplified version of the Boyer-Moore algorithms.
- Time complexity $O(mn)$.
- However, very fast in practice.

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

## The BMS Algorithm - 1

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The BMS Algorithm - 2

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The BMS Algorithm - 3

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The BMS Algorithm - 4

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The BMS Algorithm - 5

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
**The Sunday Adaption of Boyer-Moore Algorithm**
The Shift-And Algorithm
The Franek-Jennings-Smyth Algorithm

# The BMS Algorithm - 6

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
The Sunday Adaption of Boyer-Moore Algorithm
**The Shift-And Algorithm**
The Franek-Jennings-Smyth Algorithm

## The Shift-And Algorithm

- Makes use of the bit-parallel nature of computer.
- Time complexity $O(mn/w)$.
- Can be easily modified for indeterminate pattern-matching.

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
The Sunday Adaption of Boyer-Moore Algorithm
**The Shift-And Algorithm**
The Franek-Jennings-Smyth Algorithm

## The Shift-And Algorithm

| $m\backslash\Sigma$ | $A$ | $C$ | $G$ | $T$ |
|---|---|---|---|---|
| $A$ | 1 | 0 | 0 | 0 |
| $A$ | 1 | 0 | 0 | 0 |
| $T$ | 0 | 0 | 0 | 1 |
| $C$ | 0 | 1 | 0 | 0 |
| $G$ | 0 | 0 | 1 | 0 |

Preprocessing();
*BitArray*: $D[1..n]$
$D[1] \leftarrow S_{x[1]}$
**for** $i = 2$ **to** $n$ **do**
    $D[i] \leftarrow (\text{Shift}(D[i-1]) \& S_{x[i]})$;
    **if** $D_m \& 10^{m-1}$ **then** output $i - m + 1$;

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
**The Franek-Jennings-Smyth Algorithm**

## The Franek-Jennings-Smyth Algorithm

- A hybrid algorithm that combines the KMP and BMS algorithm.
- Inherits the merits of both algorithms: very fast both asymptotically ($O(n)$) and in practice.

Introduction
**Fundamental Algorithms**
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
Conclusion

The Knuth-Morris-Pratt Algorithm
The Sunday Adaption of Boyer-Moore Algorithm
The Shift-And Algorithm
**The Franek-Jennings-Smyth Algorithm**

## Outline of the FJS Algorithm

1. Perform **Sunday** shift along text.
2. When a match of letters is found at the end of the pattern, switch to **KMP** matching.
3. Continue **KMP** matching until no border can be used, then switch back to **Sunday** shift.

Introduction
Fundamental Algorithms
**Special Properties of Indeterminate Borders**
The New Hybrid Algorithm
Experimental Results
Conclusion

# Outline

Introduction
Fundamental Algorithms
**Special Properties of Indeterminate Borders**
The New Hybrid Algorithm
Experimental Results
Conclusion

## Example of Non-transitivity Effect

Suppose we are performing KMP matching along the text.

| Index | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | |
|---|---|---|---|---|---|---|---|---|---|
| $x$ | ...... | $a$ | $a$ | $b$ | $b$ | $a$ | $b$ | $b$ | ...... |
| $p$ | | $a$ | $*$ | $*$ | $b$ | $a$ | $*$ | $a$ | |
| 1$st$ Shift | | | $a$ | $*$ | $*$ | $b$ | $a$ | | ...... |
| 2$nd$ Shift | | | | | $a$ | $*$ | $*$ | | ...... |
| 3$rd$ Shift | | | | | | $a$ | $*$ | | ...... |

Table: First example of the non-transitivity effect

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

### Proposition

*Shifting the pattern to the right according to the longest border cannot guarantee a prefix match.*

### Proposition

*A border of a border of $x$ is not necessarily a border of $x$.*

Introduction
Fundamental Algorithms
**Special Properties of Indeterminate Borders**
The New Hybrid Algorithm
Experimental Results
Conclusion

| Index | | 1 | 2 | 3 | 4 | 5 | 6 | 7 | |
|---|---|---|---|---|---|---|---|---|---|
| $x$ | ...... | $a$ | $b$ | $a$ | $*$ | $a$ | $*$ | $a$ | ...... |
| $p$ | | $a$ | $b$ | $a$ | $a$ | $a$ | $b$ | $b$ | |
| Wrong Shift | | | | | | $a$ | $b$ | $a$ | ...... |
| Correct Shift | | | | $a$ | $b$ | $a$ | $a$ | $a$ | ...... |

Table: Second example of the non-transitivity effect

### Proposition

*Shifting the pattern to the right according to the longest border can miss some occurrences in between.*

Introduction
Fundamental Algorithms
**Special Properties of Indeterminate Borders**
The New Hybrid Algorithm
Experimental Results
Conclusion

## Impact

- Because of these properties, transforming regular pattern-matching algorithms that use border arrays into indeterminate pattern-matching algorithms is non-trivial (KMP, FJS, etc.)

- However, since some of these regular algorithms are very fast in practice and have nice properties, we are motivated to invent indeterminate versions of them that avoid using border arrays.

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
Examples

# Outline

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

**Outline of the New Algorithm**
Shift-And Matching
Sunday-Shift
Examples

## A New Hybrid Algorithm

We propose a new hybrid algorithm that uses Shift-And and BMS as complementary shift engines.

- 1. Perform **Sunday shift** along text.
- 2. When a match of letters is found at the end of the pattern, switch to **Shift-And matching**.
- 3. Continue **Shift-And matching** until no match can be found at the current position ($D = 0$), then skip to next possible position and switch back to **Sunday** shift.

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Outline of the New Algorithm**
**Shift-And Matching**
**Sunday-Shift**
**Examples**

## Shift-And Preprocessing

The usual Shift-And preprocessing is modified as follows:

**for** $i = 1$ **to** m
  **for** $j = 1$ **to** $|\Sigma|$
    **if** MATCH($p[i], \Sigma[j]$) **then** $S[i,j] = 1$
    **else** $S[i,j] = 0$

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

Outline of the New Algorithm
**Shift-And Matching**
Sunday-Shift
Examples

## Properties of Shift-Or

Notice some of the important properties of Shift-Or.

### Proposition

$D[j] = 1 \Leftrightarrow p[1..j] \approx x[i - j + 1..i]$

### Proposition

$D = 0$ *if and only if there doesn't exist any* $j \in 1..m$ *such that* $p[1..j] \approx x[i - j + 1..i]$

These properties enables us to move the pattern beyond $x[i]$ when we finish `ShiftAnd-Match`.

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
**Shift-And Matching**
Sunday-Shift
Examples

## ShiftAnd-MATCH

$D \leftarrow 0$
**do**
   $D \leftarrow (D \ll 1) \& S_{x[i]}$
   **if** $D \& 10^m \neq 0$ **then output** $i$
   $i \leftarrow i + 1$
*//If D =0, terminate loop according to previous proposition*
**while** ($i \leq n$ **and** $D \neq 0$)

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Outline of the New Algorithm**
**Shift-And Matching**
**Sunday-Shift**
**Examples**

## BMS Preprocessing

The usual BMS preprocessing is modified as follows.

**for** $i = 1$ **to** $|\Delta|$
   $\Delta[i] = m + 1$
**for** $i = 1$ **to** $m$
     **for** $j = 1$ **to** $|\Sigma|$
       **if** MATCH($p[i], \Sigma[j]$) **then** $\Delta[p[i]] = i$

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Outline of the New Algorithm**
**Shift-And Matching**
**Sunday-Shift**
**Examples**

## Sunday-Shift

**while not** $\text{MATCH}(p[m], x[i'])$ **do**
$\quad i' \leftarrow i' + \Delta \big[x[i'+1]\big]$
$\quad$ **if** $i' > n$ **then return**

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
**Sunday-Shift**
Examples

## Algorithm Shift-And/Sunday

$i' \leftarrow m; m' \leftarrow m - 1;$
**while** $i' \leq n$ **do**
   Sunday-Shift();
   $i \leftarrow i' - m';$
   *//After Sunday-Shift stops, perform* ShiftAnd-MATCH
   ShiftAnd-Match();
   *//After ShiftAnd-Match stops, shift pattern to the right*
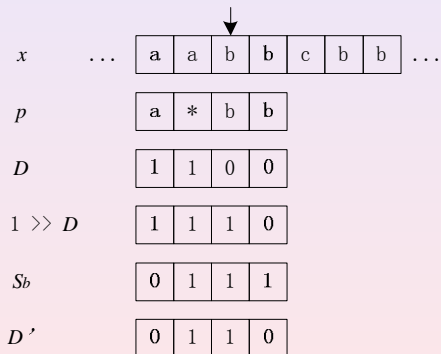   $i' \leftarrow i + m';$

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm



| $x$ | ... | a | a | b | **b** | c | b | b | ... |
|---|---|---|---|---|---|---|---|---|---|

| $p$ | | **a** | * | b | **b** |
|---|---|---|---|---|---|

| $D$ | | 1 | 0 | 0 | 0 |
|---|---|---|---|---|---|

| $1 >> D$ | | 1 | 1 | 0 | 0 |
|---|---|---|---|---|---|

| $S_a$ | | 1 | 1 | 0 | 0 |
|---|---|---|---|---|---|

| $D'$ | | 1 | 1 | 0 | 0 |
|---|---|---|---|---|---|

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm



| $x$ | $\ldots$ | **a** | a | b | **b** | c | b | b | $\ldots$ |
|---|---|---|---|---|---|---|---|---|---|

| $p$ | | **a** | * | b | **b** |
|---|---|---|---|---|---|

| $D$ | | 1 | 1 | 0 | 0 |
|---|---|---|---|---|---|

| $1 \gg D$ | | 1 | 1 | 1 | 0 |
|---|---|---|---|---|---|

| $S_b$ | | 0 | 1 | 1 | 1 |
|---|---|---|---|---|---|

| $D\,'$ | | 0 | 1 | 1 | 0 |
|---|---|---|---|---|---|

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm



| $x$ | $\ldots$ | a | a | b | **b** | c | b | b | $\ldots$ |
|-----|----------|---|---|---|-------|---|---|---|----------|

| $p$ | | a | * | b | **b** |
|-----|---|---|---|---|-------|

| $D$ | | 0 | 1 | 1 | 0 |
|-----|---|---|---|---|---|

| $1 \gg D$ | | 1 | 0 | 1 | 1 |
|-----------|---|---|---|---|---|

| $S_b$ | | 0 | 1 | 1 | 1 |
|-------|---|---|---|---|---|

| $D'$ | | 0 | 0 | 1 | 1 |
|------|---|---|---|---|---|

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm



| | | | | | | | |
|---|---|---|---|---|---|---|---|
| $x$ ... | **a** | a | b | **b** | c | b | b | ... |

| $p$ | **a** | * | b | **b** |
|---|---|---|---|---|

| $D$ | 0 | 0 | 1 | 1 |
|---|---|---|---|---|

| $1 \gg D$ | 1 | 0 | 0 | 1 |
|---|---|---|---|---|

| $S_c$ | 0 | 1 | 0 | 0 |
|---|---|---|---|---|

| $D'$ | 0 | 0 | 0 | 0 |
|---|---|---|---|---|

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Outline of the New Algorithm**
**Shift-And Matching**
**Sunday-Shift**
**Examples**

# Example of The New Hybrid Algorithm



Begin *Sunday-Shift*

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
**The New Hybrid Algorithm**
Experimental Results
Conclusion

Outline of the New Algorithm
Shift-And Matching
Sunday-Shift
**Examples**

# Example of The New Hybrid Algorithm

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

# Outline

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
**Experimental Results**
Conclusion

- In all of these tests, the hybrid algorithm's behaviour is very close to that of the better of BMS and Shift-And.

- The new algorithms's total running time is very competitive among these three algorithms being tested.

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

# Outline

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

- A new algorithm that performs fast pattern-matching on both regular and indeterminate strings.

- Strong ability to adapt to the nature of text/pattern and to achieve faster performance over cases that arise in practice. This dynamic adaptivity is useful when we do not know the type of text or pattern: we don't need to make a decision ahead of time about which algorithm to use.

- Future work: Indeterminate pattern-matching algorithms based on variants of Shift-And such as BNDM and [Fre07], as well as on new convolution techniques [AAR07].

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
**Conclusion**

📕 Avivit Levy Amihood Amir and Liron Reuveni.
The practical efficiency of convolutions in pattern matching.
*Fundamenta Informatica*, page to appear, 2007.

📕 Robert S. Boyer and J. S Strother Moore.
A fast string searching algorithm.
*CACM*, 20(10):762–772, 1977.

📕 R.A. Baeza-Yates and G.H. Gonnet.
A new approach to text searching.
*Communications of the ACM*, 35(10):74–82, 1992.

📕 Bálint Dömölki.
A universal computer system based on production rules.
*BIT*, 8:262–275, 1968.

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
**Conclusion**

📕 Frantisek Franek, Christopher G. Jennings, and W. F. Smyth.
A simple fast hybrid pattern-matching algorithm.
*Journal of Discrete Algorithms*, to appear, 2006.

📕 Kimmo Fredriksson.
Linear worst case time bndm.
*Information Processing Letters*, page to appear, 2007.

📕 Jan Holub, W. F. Smyth, and Shu Wang.
Hybrid pattern-matching algorithms on indeterminate strings.
*London Algorithmics and Stringology 2006, J. Daykin, M. Mohamed and K. Steinhoefel (eds.), King's College London Series Texts in Algorithmics*, pages 115–133, 2006.

**Introduction**
**Fundamental Algorithms**
**Special Properties of Indeterminate Borders**
**The New Hybrid Algorithm**
**Experimental Results**
**Conclusion**

📕 Jan Holub, W. F. Smyth, and Shu Wang.
Fast pattern-matching on indeterminate strings.
*Journal of Discrete Algorithms*, to appear, 2006.

📕 D. E. Knuth, J. H. Morris, and V.R. Pratt.
Fast pattern matching in strings.
*SIAM Journal on Computing*, 6(2):323–350, 1977.

📕 G. Navarro and M. Raffinot.
A bit-parallel approach to suffix automata: Fast extended
string matching.
In M. Farach-Colton, editor, *Proceedings of the 9th Annual
Symposium on Combinatorial Pattern Matching*, number
1448, pages 14–33, Piscataway, NJ, 1998. Springer-Verlag,
Berlin.

Introduction
Fundamental Algorithms
Special Properties of Indeterminate Borders
The New Hybrid Algorithm
Experimental Results
**Conclusion**

S. Wu and U. Manber.
Fast text searching with errors.
*Communications of the ACM*, 35(10):83–91, 1992.