



**Claudio Gomes de Mello**

**Codificação livre de prefixo para  
cripto-compressão**

**Tese de Doutorado**

Tese apresentada ao Programa de Pós-graduação em Informática do Departamento de Informática da PUC-Rio como requisito parcial para obtenção Do título de Doutor em Informática

Orientador: Prof. Ruy Luiz Milidiú

Rio de Janeiro  
Setembro de 2006



**Claudio Gomes de Mello**

**Codificação livre de prefixo para  
cripto-compressão**

Tese apresentada ao Programa de Pós-graduação em Informática do Departamento de Informática do Centro Técnico Científico da PUC-Rio como requisito parcial para obtenção Do título de Doutor em Informática. Aprovada pela Comissão Examinadora abaixo assinada.

**Prof. Ruy Luiz Milidiú**

Orientador

Departamento de Informática — PUC-Rio

**Prof. Weiler Alves Finamore**

CETUC (Centro de Estudos em Telecomunicações da PUC-Rio)  
— PUC-Rio

**Prof. Marco Antonio Casanova**

Departamento de Informática — PUC-Rio

**Prof. Ricardo Dahab**

Instituto de Computação — Unicamp

**Prof. José Antônio Moreira Xexéo**

Seção de Engenharia de Computação e Telemática — IME

**Prof. José Engenio Leal**

Coordenador Setorial do Centro Técnico Científico — PUC-Rio

Rio de Janeiro, 18 de Setembro de 2006

Todos os direitos reservados. É proibida a reprodução total ou parcial do trabalho sem autorização da universidade, do autor e do orientador.

**Claudio Gomes de Mello**

Ficha Catalográfica

Mello, Claudio Gomes de

Codificação livre de prefixo para cripto-compressão / Claudio Gomes de Mello; orientador: Ruy Luiz Milidiú. — Rio de Janeiro : PUC–Rio, Departamento de Informática, 2006.

v., 110 f: il. ; 29,7 cm

1. Tese (doutorado) - Pontifícia Universidade Católica do Rio de Janeiro, Departamento de Informática.

Inclui referências bibliográficas.

1. Informática – Tese. 2. Segurança da Informação. 3. Compressão de Dados. 4. Códigos de Prefixo. 5. Substituição Homofônica. I. Milidiú, Ruy Luiz. II. Pontifícia Universidade Católica do Rio de Janeiro. Departamento de Informática. III. Título.

CDD: 004

## Agradecimentos

Ao meu orientador Prof. Ruy Luiz Milidiú por todo o apoio, motivação e sabedoria. O Prof. Ruy soube guiar esse trabalho de pesquisa desde o seu início e foi essencial naqueles momentos em que o trabalho apresentava obstáculos para sua continuidade. Agradeço ao Prof. Ruy por ter passado sua experiência como pesquisador e por sua dedicação.

À minha esposa Dane Schemidt de Mello por seu apoio, sem o qual um trabalho de tanto tempo não teria como ser concluído.

Aos meus companheiros Zé Rodrigues e Carlos Eduardo que desenvolveram duas dissertações de mestrado na mesma área, contribuindo para a troca de experiências e para a evolução do trabalho.

Aos meus colegas da PUC-Rio: Frederico, Artur, Flavio, Raul, Lorenza, Duarte entre outros, que criaram um ótimo ambiente de trabalho.

Ao Exército Brasileiro pela oportunidade de cursar o doutorado, sem a qual este trabalho não existiria.

À Seção de Engenharia de Computação e Telemática (SE/8) do Instituto Militar de Engenharia (IME), pela compreensão e apoio, cedendo inclusive algumas horas de trabalho para que eu pudesse desenvolver esse trabalho.

Aos meus companheiros de IME pelo apoio total na reta final, em especial Salles, Luiz Henrique, Cecilio, Carla e Freire.

## Resumo

Mello, Claudio Gomes de; Milidiú, Ruy Luiz. **Codificação livre de prefixo para cripto-compressão**. Rio de Janeiro, 2006. 110p. Tese de Doutorado — Departamento de Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Cifragem e compressão de dados são funcionalidades essenciais quando dados digitais são armazenados ou transmitidos através de canais inseguros. Geralmente, duas operações sequenciais são aplicadas: primeiro, compressão de dados para economizar espaço de armazenamento e reduzir custos de transmissão, segundo, cifragem de dados para prover confidencialidade. Essa solução funciona bem para a maioria das aplicações, mas é necessário executar duas operações caras, e para acessar os dados, é necessário primeiro decifrar e depois descomprimir todo o texto cifrado para recuperar a informação. Neste trabalho são propostos algoritmos que realizam tanto compressão como cifragem de dados. A primeira contribuição desta tese é o algoritmo ADDNULLS — Inserção Seletiva de Nulos. Este algoritmo usa a técnica da esteganografia para esconder os símbolos codificados em símbolos falsos. É baseado na inserção seletiva de um número variável de símbolos nulos após os símbolos codificados. É mostrado que as perdas nas taxas de compressão são relativamente pequenas. A segunda contribuição desta tese é o algoritmo HHC — Huffman Homofônico-Canônico. Este algoritmo cria uma nova árvore homofônica baseada na árvore de Huffman canônica original para o texto de entrada. Os resultados dos experimentos são mostrados. A terceira contribuição desta tese é o algoritmo RHUFF — Huffman Randomizado. Este algoritmo é uma variante do algoritmo de Huffman que define um procedimento de cripto-compressão que aleatoriza a saída. O objetivo é gerar textos cifrados aleatórios como saída para obscurecer as redundâncias do texto original (confusão). O algoritmo possui uma função de permutação inicial, que dissipa a redundância do texto original pelo texto cifrado (difusão). A quarta contribuição desta tese é o algoritmo HSPC2 — Códigos de Prefixo baseados em Substituição Homofônica com 2 homofônicos. No processo de codificação, o algoritmo adiciona um bit de sufixo em alguns códigos. Uma chave secreta e uma taxa de homofônicos são parâmetros que controlam essa inserção. É mostrado que a quebra do HSPC2 é um problema NP-Completo.

### Palavras-chave

Segurança da Informação. Compressão de Dados. Códigos de Prefixo. Substituição Homofônica.

## Abstract

Mello, Claudio Gomes de; Milidiú, Ruy Luiz. **Crypto-compression prefix coding**. Rio de Janeiro, 2006. 110p. PhD Thesis — Department of Informática, Pontifícia Universidade Católica do Rio de Janeiro.

Data compression and encryption are essential features when digital data is stored or transmitted over insecure channels. Usually, we apply two sequential operations: first, we apply data compression to save disk space and to reduce transmission costs, and second, data encryption to provide confidentiality. This solution works fine for most applications, but we have to execute two expensive operations, and if we want to access data, we must first decipher and then decompress the ciphertext to restore information. In this work we propose algorithms that achieve both compressed and encrypted data. The first contribution of this thesis is the algorithm ADDNULLS — Selective Addition of Nulls. This algorithm uses steganographic technique to hide the real symbols of the encoded text within fake ones. It is based on selective insertion of a variable number of null symbols after the real ones. It is shown that coding and decoding rates loss are small. The disadvantage is ciphertext expansion. The second contribution of this thesis is the algorithm HHC — Homophonic-Canonic Huffman. This algorithm creates a new homophonic tree based upon the original canonical Huffman tree for the input text. It is shown the results of the experiments. Adding security has not significantly decreased performance. The third contribution of this thesis is the algorithm RHUFF — Randomized Huffman. This algorithm is a variant of Huffman codes that defines a crypto-compression algorithm that randomizes output. The goal is to generate random ciphertexts as output to obscure the redundancies in the plaintext (confusion). The algorithm uses homophonic substitution, canonical Huffman codes and a secret key for ciphering. The secret key is based on an initial permutation function, which dissipates the redundancy of the plaintext over the ciphertext (diffusion). The fourth contribution of this thesis is the algorithm HSPC2 — Homophonic Substitution Prefix Codes with 2 homophones. It is proposed a provably secure algorithm by using a homophonic substitution algorithm and a key. In the encoding process, the HSPC2 function appends a one bit suffix to some codes. A secret key and a homophonic rate parameters control this appending. It is shown that breaking HSPC2 is an NP-Complete problem.

## Keywords

Information Security. Data Compression. Prefix Codes. Homophonic Substitution.

# Sumário

1	Introdução	12
2	Teoria da Informação	14
2.1	Alfabeto, texto, letras e caracteres	14
2.2	Partição e cadeia	14
2.3	Fonte de Informação	14
2.4	Entropia	14
2.5	Redundância	16
2.6	Códigos	17
2.7	Distância de Unicidade	20
3	Sistemas Criptográficos	22
3.1	Introdução	22
3.2	Conceitos iniciais	22
3.3	Sistemas Criptográficos Clássicos	27
3.4	Substituição Homofônica	30
3.5	Sistemas Criptográficos Simétricos	33
3.6	Sistemas Criptográficos de Chave Pública	33
3.7	Geração de chaves	34
3.8	Funções de Resumo da Mensagem	34
3.9	Assinaturas Digitais	35
4	Códigos de Huffman	37
4.1	Árvores de código	37
4.2	Códigos de Huffman	38
4.3	Códigos de Huffman canônicos	43
5	Algoritmos Cripto-Compressores	45
5.1	Introdução	45
5.2	Trabalhos Correlatos	45
5.3	Algoritmos Propostos	46
6	Inserção Seletiva de Nulos	48
6.1	Introdução	48
6.2	Huffman Criptografado	48
6.3	Experimentos	55
6.4	Conclusões	56
7	Huffman Homofônico-Canônico	57
7.1	Introdução	57
7.2	Esquema geral de funcionamento	57
7.3	O Processo de Codificação do algoritmo HHC	58
7.4	O Processo de Decodificação do HHC	64
7.5	O Esquema de Utilização do Algoritmo	65
7.6	Algoritmo Huffman Canônico	69



7.7	Resultados Experimentais	69
8	Huffman Randomizado	<b>82</b>
8.1	Introdução	82
8.2	Descrição do Algoritmo	82
8.3	Experimentos	86
8.4	Conclusões	89
9	Códigos de Prefixo baseados em Substituição Homofônica com 2 homofônicos	<b>90</b>
9.1	Introdução	90
9.2	Função de Substituição Homofônica	90
9.3	Quebrando o Código	95
9.4	Estimando a Expansão do Texto Cifrado	102
9.5	Conclusões	103
10	Conclusão e Trabalhos Futuros	<b>105</b>

## Lista de figuras

3.1	Um esquema geral para substituição homofônica.	30
3.2	Substituição homofônica convencional.	31
3.3	Substituição homofônica de tamanho variável.	31
3.4	Esquema de autenticação digital.	36
4.1	Representação de códigos binários por árvores binárias.	38
4.2	(a) Passo 1 - Início da construção da árvore de Huffman.	39
4.3	(b) Passo 2.	40
4.4	(c) Passo 3.	40
4.5	(d) Passo 4.	40
4.6	(e) Passo 5 - Árvore de Huffman final.	40
4.7	Navegação na árvore de Huffman.	41
4.8	Árvore de Huffman canônica.	44
6.1	Geração de uma nova árvore de Huffman falsa a partir da taxa $\alpha$ .	50
6.2	Geração dos homofônicos de $\eta$ seguidos pelo símbolo efetivo.	51
6.3	XOR entre o $i$ -ésimo bit da chave $k$ e o ID bit da palavra código.	54
6.4	Chave composta.	55
7.1	Visão Geral do Algoritmo <i>HHC - Huffman Homofônico Canônico</i> .	58
7.2	Processo de Codificação do <i>HHC</i> .	59
7.3	Árvore canônica para 3 símbolos homofônicos de probabilidades iguais.	61
7.4	Permutação no nível da árvore.	63
7.5	Processo de Decodificação.	64
7.6	Esquema de Codificação <i>Tradicional</i> .	66
7.7	Esquema de Codificação <i>Modelo Codificado Vulnerável</i> .	67
7.8	Esquema de Codificação <i>Modelo Codificado Protegido</i> .	67
7.9	Esquema de Codificação <i>Chave Composta com o Modelo</i> .	68
7.10	Velocidade de Codificação do arquivo <i>alencar.txt</i> .	72
7.11	Aumento do tamanho do tamanho do <i>Arquivo Codificado e Arquivo Modelo</i> para o arquivo <i>alencar.txt</i> .	72
7.12	Perda de Probabilidade do arquivo <i>alencar.txt</i> .	73
7.13	Comparação entre as taxas de compressão do <i>HHC</i> e do <i>Huffman Canônico</i> .	73
7.14	Comparação entre as taxas de compressão do <i>Homofônico Canônico</i> e do <i>Huffman Canônico</i> (Fase 1 e 2).	75
8.1	Arquitetura do algoritmo	83

## Lista de tabelas

2.1	Distribuição de probabilidades típica (língua inglesa).	17
2.2	Exemplo de códigos binários.	18
4.1	Frequências por símbolo.	39
4.2	Códigos de Huffman por símbolo.	41
4.3	Códigos de Huffman canônico.	44
6.1	Resultados da expansão de texto(tamanhos em bytes).	56
7.1	Descrição dos <i>arquivos de teste</i> .	70
7.4	Simulação do número de homofônicos com o arquivo <i>alencar.txt..</i>	76
7.6	Codificação do <i>Arquivo Original</i> com o <i>Huffman Canônico</i> .	77
7.8	Codificação do <i>Arquivo Original</i> com o <i>HHC</i> .	77
7.10	Tempos de Codificação e Decodificação.	78
7.12	Tempos de Codificação por módulo.	78
7.14	Codificação do <i>Arquivo Modelo</i> usando <i>Huffman Canônico</i> .	79
7.16	Codificação do <i>Arquivo Modelo</i> usando o <i>HHC</i> .	79
7.18	Tempo de Codificação do <i>Arquivo Modelo</i> .	80
7.20	Final da Codificação usando <i>Huffman Canônico</i> .	80
7.22	Final da Codificação usando o <i>HHC</i> .	81
7.24	Tempos de Codificação e Decodificação Finais.	81
8.1	Frequências relativas a $T = bcbaabaacaabaababcaacac$ .	85
8.2	Arquivos texto da coleção de GUTENBERG.	87
8.3	Arquivos texto em detalhe.	88
8.4	Número de símbolos/bloco.	88
8.5	Comparando o desempenho bits/palavra.	88
8.6	Expansão de símbolos nulos.	88
9.1	Livro de códigos original.	91
9.2	Livro de códigos homofônicos.	91
9.3	Transformação Homofônica.	92
9.4	Taxas Homofônicas.	92