

# 유니버살 데이터 압축의 최근의 연구동향

正會員 朴 志 煥\* 正會員 陳 庸 玉\*\*

## Recent Trends of Universal Data Compression

Ji Hwan PARK\*, Yong Ohk CHIN\*\* *Regular Members*

**要 約** 텍스트 데이터 중심의 컴퓨터 통신의 활성화와 멀티 미디어의 등장에 따른 영상 데이터의 취급은 대용량의 기억공간과 전송시간을 요구하게 되어 데이터 압축의 중요성이 더욱 크게 되었다. 정보원의 사전 지식을 전제로 하지 않고도 효율적인 압축을 달성하는 유니버살 부호는 동일 알고리즘으로 다양한 데이터에 적용할 수 있는 만능성을 갖는 적응성이 뛰어난 방식으로 최근 주목되고 있다.

이 논문에서는 유니버살 부호의 기본원리와 분류를 통한 실험 알고리즘에 대하여 알아보고 그 응용을 중심으로 한 최근의 연구활동에 관하여 기술한다.

**ABSTRACT** Data compression has important application in the areas of file storage and distributed computer systems. The universal data compression achieves asymptotically optimum compression ratio for strings generated by any stationary ergodic source without a priori source probabilities.

This paper describes the principle and the recent research trends on universal data compression and its applications.

### I. 서 론

정보화 사회의 발달에 따라 컴퓨터에서 처리하는 데이터는 대량화 및 다양화 되어 가고 있다. 이와같은 각종 데이터를 효율 좋게 저장하고 전송하기 위한 방법으로 데이터 압축 기술이 여러가지 고안되어 실용화 되고 있다.

데이터 압축은 원래의 데이터 계열에 포함되어 있는 불필요한 부분(redundancy)을 제거하여 유용한 정보만을 추출함으로써 보다 짧은 데이터 계열로 변환하는 것이다. 압축 데이터 계열로부터 원래의 데이터 계열로 재생할때 복원 오류를 허용하지 않는 무잡음 압축(noiseless compression) 과 일정 허용 범위내에서 오류를 허용하는 유잡음 압축(compression with distortion)으로 분류된다. 유잡음 압축은 연속 정보인 TV,

Radio 및 Sonar 신호 등의 1차 압축에 응용되며 rate distortion 함수의 해석이론에 그 이론적 바탕을 두고 있다.<sup>(1)</sup> 무잡음 압축은 데이터 처리 계나 디지털 백시밀리등의 이산 정보를 대상으로 하며 Shannon의 정보론에 그 기초를 두고 있다.<sup>(2)</sup> 이 논문에서는 이산 정보를 대상으로 하는 무잡음 압축에 한하여 다루기로 한다.

Huffman부호<sup>(3)</sup>를 비롯한 기존의 많은 압축법에서는 정보원 기호(source symbols)의 출현빈도를 알고 있다는 전제하에 최적 부호(optimum code)를 구성하는 통계적 기법을 사용하고 있다. 그러나 현실의 정보원으로부터 발생하는 많은 데이터는 그 확률 구조를 알지 못하며 심지어 마르코프 모델로 표현할 수 있는지 없는지조차 알지 못하는 경우가 대부분이다. 이와같은 배경아래 등장된 유니버살 데이터 압축법은 동일한 알고리즘으로 보다 넓은 클래스의 정보원에 적용 가능한 만능성을 가진다.<sup>(4)</sup> 또한 데이터의 길이가 충분히 긴 경우, 평균 부호의 길이를

\*釜山水産大學校 電子計算學科  
Natl. Fisheries Univ. of PUSAN

\*\*慶熙大學校 電子工學科  
Kyunghee University  
論文番號 : 91-84

압축의 이론적 한계인 엔트로피에 얼마든지 접근시킬 수 있는 점근적 최량성(asymptotically optimal) 이 보장되는 강력한 부호로 주목되고 있다.

유니버살 부호의 명확한 정의는 확립되어 있지 않으나, 첫째 어떤 클래스에 속하는 모든 정보원 출력을 최적 또는 효율 좋게 압축할 수 있는 부호 둘째 정보원의 확률분포에 의존하지 않는 부호화·복호화 알고리즘을 이용하고 있는 부호들의 의미를 포함하고 있으며 크게 다음의 두 가지로 분류된다.<sup>(6)</sup>

- (1) 데이터 계열로부터 정보원 모델링(source modeling)을 적응적으로 수행하면서 그 모델에 기초한 확률분포를 이용하여 부호화하는 방법
- (2) 현재의 부분계열이 처리가 끝난 이전의 데이터 계열의 어느 부분과 일치하고 있는가를 나타내는 정보를 이용하여 알고리즘적으로 부호화하는 방법

(1)의 구체적인 부호화를 위하여 산술부호(arithmetic codes)를 이용하는 방법<sup>(6a)</sup>과 Dynamic Huffman부호를 이용하는 방법<sup>(6b)</sup>이 있다. (2)에는 Ziv-Lemple부호<sup>(8)</sup>, Recency Rank부호<sup>(9)</sup> 및 Interval 부호<sup>(10)</sup> 등이 있다. 산술부호는 주어진 문맥모델에서의 빈도 분포가 결정되면 엔트로피까지 압축이 가능하므로 관심의 초점은 어떻게 좋은 문맥모델(context model)을 작성하느냐에 집중된다. 또한 적응사전(adaptive dictionary)을 이용하는 (2)의 방법은 먼저 데이터 계열을 적응사전을 이용하여 자연수의 계열로 변환하고, 다음에 그 자연수를 2진계열로 부호화하는 2단계로 이루어진다. 문맥모델을 이용하는 (1)의 방식과 적응사전을 이용하는 (2)의 방식은 전혀 다르게 보이나 정보원 모델 구축의 입장에서 보면 서로 밀접한 관계를 갖는 것으로 간주할 수 있다. 일반적으로 산술 부호의 최적성에 의해 문맥모델을 이용하는 방식이 높은 압축 효과를 얻을 수 있으나, 부호화·복호화의 고속성이나 소요 기억용량의 면에서 적응사전을 이용하는 방식이 경제적이므로<sup>(11)</sup> 실제에는 (2)의 방식이 널리

이용되고 있는 실정이다.

이 논문에서는 위에 언급한 각종 유니버살 부호의 기본원리와 그 알고리즘을 중심으로 다루고 이에 관련된 최근의 연구동향에 대하여 기술한다.

## II. Dynamic Huffman부호

Huffman부호는 평균 부호의 길이를 최소화하는 최단부호(compact codes)로 잘 알려져 있으나 이 부호를 이용한 데이터 압축 알고리즘의 구성에는 두가지 방법이 있다.

### 2.1 이중 주사 방식(two pass method)<sup>(9)</sup>

길이  $n$ 의 데이터 계열  $X_n = x_1 x_2 x_3 \dots x_n$ 으로부터 각 기호의 출현빈도를 구하여 그 빈도분포에 기초한 Huffman tree를 작성한다. 부호화는 먼저 Huffman tree 또는 빈도표를 부호화한 후, 작성된 Huffman tree를 이용하여  $X_n$ 을 Huffman 부호화 한다.

### 2.2 단일 주사 방식(one pass method)

$X_n$ 의 한 기호를 처리할때 마다 그때까지의 빈도분포에 의해 결정되는 tree가 Huffman tree가 되도록 어느 정해진 산법으로 tree를 매회 갱신하면서 Huffman부호화를 실시한다. Dynamic Huffman부호는 Faller<sup>(12)</sup>와 Gallager<sup>(13)</sup>에 의해 독립적으로 제안되었으며, 그후 Knuth<sup>(7)</sup>와 Vitter<sup>(14)</sup>에 의해 상세한 알고리즘이 제시되었다.

Dynamic Huffman부호의 code tree 수정 및 부호화의 일예를 Kunth 알고리즘에 따라 그림1에 나타낸다. 지금 계열  $X_3 = abc$ 에 대해 그림1(a)와 같이 tree가 작성되어 있다고 하자. 각 단점(terminal nodes)에는 지금까지 출현한 기호와 출현빈도가 표시되어 있고 각 내부절점(internal nodes)에는 그 절점에 연속되는 절점의 빈도 합이 기록되어 있다. 이때 다음문자  $x_4 = c$  라면 101이 c에 대응하는 부호어(codeword)가

되고, 아직 출현하지 않은 기호  $x_4=d$ 가 나타나면  $100(d)_2$ 를 부호어로 한다. 단,  $(d)_2$ 는 기호  $d$ 를 2진수로 나타낸 것이다. Tree의 갱신을 위해 출현 기호에 대응하는 절점의 빈도 값을 1 증가시키나, 만약 증가전의 절점과 같은 빈도를 갖는 절점이 고순위(순위는 위 일수록 높고 같은 레벨이면 오른쪽일수록 높다)에 있으면 sub-tree를 서로 교환한다. 그림 1(a)의 tree를  $x_4=c$  및  $x_4=d$ 로 갱신한 후의 tree를 그림1(b),(c)에 나타낸다.

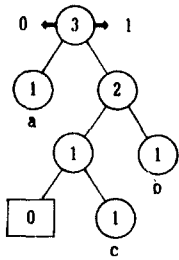


Fig 1. (a) Dynamic Huffman tree after processing "abc".

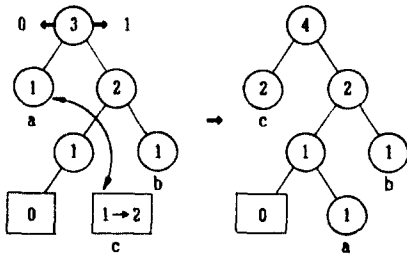


Fig 1. (b) Dynamic Huffman tree after processing "abcc".

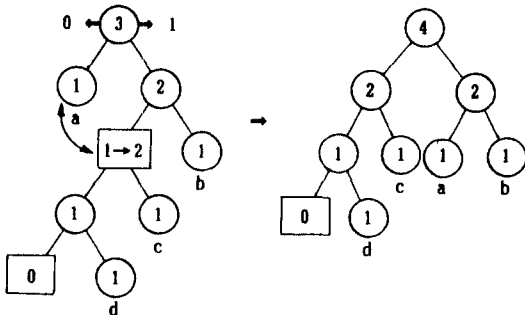


Fig 1. (c) Dynamic Huffman tree after processing "abcd".

Knuth의 알고리즘을 이용하면 tree가 항상 Huffman tree로 되기 때문에 평균 부호의 길이는 최단이 된다. 그러나 서로 교환하는 절점이 동일 빈도의 최하위 순위로 되기 때문에 최대 부호의 길이가 최단이 되지 않는다. Vitter는 Schwarz의 알고리즘<sup>(15)</sup>을 이용하여 최대 부호의 길이와 전체 부호의 길이가 최소가 되고, sub-tree의 정렬 회수를 감소시키는 알고리즘을 제시하였다.<sup>(14)</sup> Dynamic Huffman 부호를 이용한 응용 가운데 Knuth의 알고리즘에 기초한 UNIX의 COMPACT command<sup>(16)</sup>등 압축용으로서 널리 이용되고 있으며, 한글 데이터의 표현<sup>(17)</sup> 및 데이터 압축과 암호화를 결합시키는 등의 연구<sup>(18)</sup>가 이루어지고 있다.

### III. 자연수의 유니버살 표현

적응사전을 이용하는 Dictionary법<sup>(19)</sup>은 현재의 데이터 계열이 처리가 끝난 계열의 어느 부분에 일치하고 있는가를 나타내는 인덱스(자연수)  $i \in \mathbb{N}^+$ 를 부호어로서 이용한다. 따라서 자연수  $i$ 를 얼마나 효율 좋게 2진 표현 하느냐가 부호화 성능에 큰 영향을 미치게 된다. 아래의 자연수의 2진 표현법은 모두 부호어 사이의 구별이 자동적으로 이루어지는 Comma free 조건을 만족하게 된다. 이하,  $(i)_2$ 는 최상위 비트가 1인 2진수이며,  $(i)^p$ 는 p비트로 나타낸  $i$ 의 2진수 표현이다. 또한 각 표현법의  $i$ 의 2진 표현을  $B(i)$ 로 나타내며 이때 필요한 비트수를  $L(i)$ 로 표시한다. 이하  $\log$ 의 밑은 2이다.

#### 3.1 자연수 $i$ 가 유한일때( $1 \leq i \leq N$ )

(1)  $i$ 가 같은 확률로 발생하는 경우

$\lceil \log N \rceil$  비트의 고정 길이로 간단히 표현(이하,  $B_1$ 법)할 수 있으나, code-tree가 complete tree로 되어 있지 않기 때문에  $(\lceil \log N \rceil - \log N)$  비트의 손실이 발생하게 된다. 여기서  $\lceil x \rceil$ 는  $x$ 보다 큰 최소 정수이며,  $\lfloor x \rfloor$ 는  $x$ 보다 작은 최대 정수이다. 이 결점을 개량한 방법<sup>(20)</sup>은 압축율의

개선 효과가 상당히 크므로 처리 속도가 중요시 되지 않는 경우에는 유효한 방법으로 기대된다. 이때  $B_1(i;N)=(i-1)^{\lceil \log N \rceil}$ 으로 표기된다.

(2)  $i$ 값이 작을수록 높은 확률을 갖는 경우 :  $P(i-1) \geq p(i)$

$i$ 의 값이 작을수록 짧은 비트로 표현할 필요가 있는데, Willems<sup>(40)</sup>는  $N$ 개의 자연수를 몇개의 그룹으로 나누어  $i$ 를  $G$ 번째 그룹의  $M$ 번째의 멤버로 나타내어  $G, M$ 를 2진 부호화하는 방식을 제안하였다. Willems부호의 표현방법(이하,  $B_2$ 법)과 그 일예를 다음에 보인다.

$$2^G \leq i \leq 2^{G+1} - 1, G=0,1,\dots,L=\lceil \log N \rceil \text{에 대해,}$$

$$B_2(i;N) = \begin{cases} (G)_{\lceil \log(i/2^G) \rceil} \cdot (M)^G, M=i-2^G, i \neq N \\ (G)_{\lceil \log(i/2^G) \rceil}, i=N \end{cases}$$

이므로 소요되는 비트수  $L_2(i) = \lceil \log(\lceil \log N \rceil + 1) \rceil + \lceil \log i \rceil, i \neq N$ 으로 된다.  $N=8$ 의 경우,

$$\begin{aligned} B_2(1;8) &= 00, & G=0 & & B_2(5;8) &= 10 \ 01, & G=2 \\ B_2(2;8) &= 01 \ 0, & G=1 & & B_2(6;8) &= 10 \ 10, & G=2 \\ B_2(3;8) &= 01 \ 1, & G=1 & & B_2(7;8) &= 10 \ 11, & G=2 \\ B_2(4;8) &= 10 \ 00, & G=2 & & B_2(8;8) &= 11, & G=3 \end{aligned}$$

로 되어  $i$ 값이 작을수록 짧은 길이로 표현할 수 있게 된다.  $-$ 는 그룹 $G$ 와 멤버 $M$ 을 식별하기 위한 가상의 기호이다.

### 3.2 $i$ 가 유한이 아니고, $i$ 값이 작을수록 높은 확률을 갖을때

#### (1) Zero Run법<sup>(21)</sup>

이 방법은 자연수  $i \in \{1,2,\dots\}$ 를 comma free 조건을 만족하도록  $(i)_2$ 의 길이  $\lceil \log i \rceil + 1$ 를 나타내기 위해  $\lceil \log i \rceil$  비트의 Zero Run을  $(i)_2$ 의 선두에 붙이는 표현법(이하,  $B_3$ 법)이다. 따라서,  $B_3(i) = 0 \dots 0 \cdot (i)_2, B_3(9) = 000 \cdot 1001$ 이므로  $L_3(i) = 1 + 2\lceil \log i \rceil$ 로 된다.

#### (2) Doubly Zero Run법<sup>(21)</sup>

Zero Run법의 0런의 길이에 대해 또다시  $Z$

ero Run법을 적용시키는 방법(이하,  $B_4$ 법)으로서  $i$ 의 값이 클때 효과적이다. 즉,  $B_4(i) = 0 \dots 0 \cdot (1 + \lceil \log i \rceil)_2 \cdot (i)_2$ 로 되며  $(i)_2$ 의 MSB는 항상 1로서 생략 가능하므로  $L_4(i) = 1 + \lceil \log i \rceil + 2\lceil \log(1 + \log i) \rceil$ 로 된다.  $B_3(i)$ 는 유니버살 부호로 최적성이 보장되지 않으나,  $B_4(i)$ 는 점근적으로 최적이 되는 부호이다. 한편,  $(i)_2$ 의 길이를 나타내는 부분을 재귀적으로 정의한  $B_5(i)$ 는 다음과 같다.

$$B_{5L}(i) = R_L(i) \cdot 0, B_{52}(20) = \underline{10} \ \underline{100} \ \underline{10100} \cdot 0,$$

$$B_{53}(20) = \underline{100} \ \underline{10100} \cdot 0.$$

따라서  $L_5 \approx 1 + \log i + \log \log i + \dots$ 로 된다.  $B_5$ 의 밑줄 친 각 블록의 MSB가 반드시 1임을 이용하여 최후에 0을 붙여 부호의 끝을 나타내고 있다. 그러나, 반복되는 부분의 비트를 반전시킨  $B'_{52}(20) = \underline{01} \ \underline{011} \ \underline{10100}$ 로 되므로 최후의 0을 생략할 수 있어 한 부호당 1 비트 절약할 수 있다.<sup>(22)</sup>

#### (3) Fibonacci code<sup>(23)</sup>

Elias에 의한 표현 방법에서는 부호어를 서로 구별하기 위해  $(i)_2$ 의 길이를 나타내는 prefix를 이용하였으나, Fibonacci code에서는 부호간의 구별을 위한 flag를 suffix에 부가하고 있다. Fibonacci number의 차수가  $r$ 일때  $j$ 번째의 수는  $r$ 개 이전의 Fibonacci number의 합으로 나타낸다. 즉,

$$F_j = F_{j-1} + F_{j-2} + \dots + F_{j-r}, \quad j \geq 1, F_0 = \dots = F_{r-1} = 1.$$

자연수  $i$ 의 Fibonacci 표현  $R(i)$ 는

$$R(i) = \sum_{j=0}^k d_j \cdot F_j, \quad d_j \in \{0,1\}, \quad k < i$$

의  $d_j$ 계열로 나타낸다. 차수  $r=2$ 의 Fibonacci code는  $F(i) = \sim R(i) \cdot 1$ 로 되며  $\sim R(i)$ 는  $R(i)$ 의 비트 순서를 역순으로 나타낸 계열이다.  $r=2$ 의 Fibonacci code의 일예를 표1에 나타낸다.

표 1. Fibonacci number에 의한 표현  
Table 1. Fibonacci representations and codes.

i	R(i)	F(i)=~R(i)·1
1	1	1 <u>1</u>
2	1 0	0 1 <u>1</u>
3	1 0 0	0 0 1 <u>1</u>
4	1 0 1	1 0 1 <u>1</u>
5	1 0 0 0	0 0 0 1 <u>1</u>
6	1 0 0 1	1 0 0 1 <u>1</u>
7	1 0 1 0	0 1 0 1 <u>1</u>
8	1 0 0 0 0	0 0 0 0 1 <u>1</u>
16	1 0 0 1 0 0	0 0 1 0 0 1 <u>1</u>
32	1 0 1 0 1 0 0	0 0 1 0 1 0 1 <u>1</u>
	F <sub>6</sub> F <sub>5</sub> F <sub>4</sub> F <sub>3</sub> F <sub>2</sub> F <sub>1</sub> F <sub>0</sub>	
	21 13 8 5 3 2 1	

Fibonacci code는 prefix code를 구성하지는 않으나 최후미의 1에 의해 comma free code의 조건을 만족한다. 또한 점근적 최량성은 보장되지 않지만, 유니버설 코드로서 i(514,228에서는 Elias code에 비해 우수한 성능을 보이고 있어 유한 길이의 데이터 압축에 많이 이용되고 있다.<sup>(24)</sup>

#### IV. 적응사전(Adaptive Dictionary:AD)을 이용한 방식

사전을 이용한 데이터 압축법은 오래전부터 사용되어 왔으나, 효율이 좋은 실용적인 적응사전을 이용한 방식의 연구는 Ziv & Lempel의 논문<sup>(8)(25)(26)</sup>이 계기가 되어 활발한 연구가 이루어지게 되었다. AD방식은 데이터 계열을 자연수 계열로 변환한 후 그것을 2진 계열로 부호화하게 된다. 데이터 계열로부터 자연수 계열로의 부호화는 다음의 3단계로 이루어지게 된다.

- (1) 검색 : 압축하고자 하는 데이터 기호(열)을 사전 내에서 검색.
- (2) 부호화 : 데이터 열을 사전 내의 순위등을 이용하여 자연수의 계열로 부호화

(3) 갱신 : 부호화 / 복호화의 사전을 갱신.

Ziv-Lempel이 포인터형의 유니버설 부호(이하 ZL-77)<sup>(25)</sup>를 제안하여 데이터 압축에 적용시킨 이래 수많은 연구가에 의해 응용되어 압축 유틸리티<sup>(27)</sup>로써 널리 사용되게 되었다.

#### 4.1 ZL-77방식

Ziv-Lempel부호의 기본원리는 원 데이터 계열의 이미 출현한 기호열로부터 기호 복제를 거듭하여 새로운 기호열을 만들면서 압축계열로 치환하는 것이다. 기호열 X=xy에 대해, X[a,b]는 X의 a의 위치에서 시작되는 길이 b의 기호열의 생성을 의미한다. 즉 X[1,1]=x, X[2,2]=xy 및 X[2,4]=xyxy로 복제된다. 또한 새로운 기호의 출현에 대응하기 위하여 새로운 기호를 추가하는 규칙을 도입한다. Y=xyxyz일때, X로부터 Y를 직접은 복제할 수 없으나 X[2,2]로부터 xyz를 생성하기 위해

$$X[2,2]=xy \Rightarrow ^2xy=xy \cdot z$$

라는 생성규칙을 이용하면 복제 가능하다. 따라서, 모든 임의의 계열은 기호(열)의 생성을 반복함으로써 서로 다른 부분계열로 분해된다. 이것에 바탕을 둔 ZL-77방식의 알고리즘은 다음과 같다.

(1) 버퍼를 초기화한다.

길이 p의 참조부 버퍼 P와 길이 q의 부호화부 버퍼 Q를 이용하여 기호복제를 수행한다. 최초에는 P에 아무 데이터도 들어 있지 않는 공백의 상태이고, Q에는 원 데이터의 선두 부분이 설정되어 있다.

(2) 부분계열로 분해한다.

생성규칙을 이용하여 P로부터 Q를 부분계열로 분해한다. 즉 계열 X=xyxyz는 (λ)=<sup>x</sup>x=<sup>y</sup>x·y=<sup>z</sup>xy·z(x·xy·xyz)의 세개의 부분계열로 분해된다. i번째의 버퍼의 상태를 B<sub>i</sub>=P·Q라 하면 Q에 대해 부분계열 X<sub>i</sub>=Q(1,b<sub>i</sub>)를 아래와 같이 구한다.

$$X_1=Q(1,b_1)=Q(1,b_1-1) \cdot Q(b_1)=P(a_1,a_1+b_1-2) \cdot Q(b_1)=P[a_1,b_1-1] \cdot Q(b_1)$$

Q(1,b<sub>1</sub>-1)은 이미 부호화가 끝난 계열로부터 복제되는 부분으로서 b<sub>1</sub>-1가 0인 경우는 공계열(λ)을 복제함을 의미한다. 여기서 P(a,b)는 계열 P의 a에서 b까지의 부분계열을 나타내며 P[a,b]와 구별됨에 주의 바란다.

(3) 부분계열을 부호화한다.

부분계열 X<sub>1</sub>에 대한 부호어 C<sub>1</sub>=C<sub>11</sub>C<sub>12</sub>C<sub>13</sub>로 이루어진다. 여기서 C<sub>11</sub>=a<sub>1</sub>-1, C<sub>12</sub>=b<sub>1</sub>-1 및 C<sub>13</sub>=Q(b<sub>1</sub>)에 해당한다.

(4) 버퍼를 갱신한다.

버퍼의 내용을 부호화 처리가 이루어진 데이터의 수 만큼 좌로 이동시키고 Q버퍼에는 이동된 기호의 수만큼 새로운 데이터가 들어오게 된다. (2)에서 (4)의 절차를 원 데이터가 전부 처리될 때까지 반복한다. 아래에 ZL 77방식에 의한 부호화의 일예를 보인다.

<전제조건>

$$\begin{aligned} \text{원 데이터} & : X = \text{xyxyzyzyzx} \cdots \cdots \\ & = x \cdot xy \cdot xyz \cdot zzyzx \cdots \cdots \end{aligned}$$

버퍼 P와 Q의 크기 : p=q=8 byte

C<sub>11</sub>과 C<sub>12</sub>의 크기 : ⌈log<sub>2</sub> p⌉=3 bit, ⌈log<sub>2</sub> q⌉=3 bit

①버퍼 : B<sub>1</sub>=P · Q, P=: :: ::, Q=xyxyzyzy

$$\text{분해} : X_1=x=Q(1,0) \cdot x=P[1,0] \cdot x$$

$$\text{부호화} : C_{11}=000, C_{12}=000, C_{13}=x$$

②버퍼 : B<sub>2</sub>=P · Q, P=: :: :: :: x, Q=xyxyzyzy

$$\text{분해} : X_2=xy=Q(1,1) \cdot y=P[8,1] \cdot y$$

$$\text{부호화} : C_{21}=111, C_{22}=001, C_{23}=y$$

③버퍼 : B<sub>3</sub>=P · Q, P=: :: :: :: xxy, Q=xyzyzyzx

$$\text{분해} : X_3=xyz=Q(1,2) \cdot z=P[7,2] \cdot z$$

$$\text{부호화} : C_{31}=110, C_{32}=010, C_{33}=z$$

④버퍼 : B<sub>4</sub>=P · Q, P=: :: xxyxyzyzyzx-

$$\text{분해} : X_4=zyzyzx=Q(1,4)=P[7,4] \cdot x$$

$$\text{부호화} : C_{41}=110, C_{42}=100, C_{43}=x$$

ZL 77방식의 주된 계산량은 일치하는 최대 길이의 기호열을 P버퍼로부터 탐색하기 위한 패턴매칭(pattern matching) 작업이다. 이를 위해 KMP 방식이나 BM방식등의 실시간 패턴매칭 알고리즘<sup>(29)</sup>을 이용한 방식이 있다.<sup>(29)</sup> ZL-77방식에 관한 많은 개량 방식이 연구 및 실용화 되어 전자우편, BBS를 비롯한 PC통신 분야와 각종 데이터의 효율적인 보관을 위해 널리 사용되고 있다.<sup>(30)</sup>

#### 4.2 ZL-78방식

Ziv & Lempel은 ZL 77방식의 유니버살 부호 이외에 IL(Information Lossless)부호기에 대한 부호화 정리를 증명하기 위하여 증분분해(incremental parsing)을 이용한 유니버살 부호(이하 ZL-78방식)를 제안하고 있다.<sup>(26)</sup> ZL-78방식의 기본개념인 증분분해는 원 데이터 계열 X<sub>n</sub>=x<sub>1</sub>x<sub>2</sub>x<sub>3</sub>...x<sub>n</sub>을 X=X<sub>0</sub>X<sub>1</sub>X<sub>2</sub>...X<sub>N</sub>(N<n)으로 분해할 때, 각 부분계열 X<sub>i</sub>(0≤i≤N)는 다음 조건을 만족한다.

(1) X<sub>0</sub>=λ(길이 0의 Null string)

(2) X<sub>N</sub>을 제외한 모든 X<sub>i</sub>(1≤i<N)는 서로 다르다.

(3) X<sub>i</sub>의 최후의 문자 x<sub>i</sub>를 제거하면 일치하는 X<sub>j</sub>(0≤j<i)가 반드시 하나 존재한다.

각 X<sub>i</sub>는 (3)의 성질에 의해 X<sub>i</sub>=X<sub>j</sub> · x<sub>i</sub>를 이용하여 (j,x<sub>i</sub>)를 부호화하면 된다. ZL-78방식의 구체적인 부호화 알고리즘은 다음과 같다.

(1) 기호열을 부분계열로 증분분해한다.

$$X_i = X_j \cdot x_i, j < i$$

(2) 부분계열을 부호화한다.

$$\text{부호어 } C_i = j \cdot r + x_i, x_i < r(\text{radix}) \cdot$$

로 이루어지며, r=2일때 C<sub>i</sub>의 upper bound는 C<sub>i</sub>=2i-1로 되어 X<sub>i</sub>를 부호화하기 위해 필요한 비트수 L<sub>i</sub>는 ⌈log i⌉+1로 된다. 아래에 X<sub>9</sub>=01010100에 대한 부호화의 일예를 보인다.

(1) 증분분해

- $X_0 = \lambda$
- (1)  $X_1 = X_0 \cdot 0 \rightarrow 0$
  - (2)  $X_2 = X_1 \cdot 1 \rightarrow 01$
  - (3)  $X_3 = X_2 \cdot 0 \rightarrow 010$
  - (4)  $X_4 = X_0 \cdot 1 \rightarrow 1$
  - (5)  $X_5 = X_1 \cdot 0 \rightarrow 00$

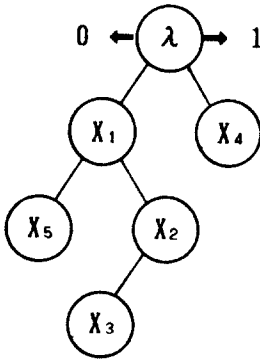


Fig 2. parsing tree

② 부호화

- (1)  $X_1 = X_0 \cdot 0 = 0, C_1 = 0 \cdot 2 + 0 = 0, L_1 = \lceil \log 1 \rceil + 1 = 1, b_1 = 0$
- (2)  $X_2 = X_1 \cdot 1 = 01, C_2 = 1 \cdot 2 + 1 = 3, L_2 = \lceil \log 2 \rceil + 1 = 2, b_2 = 11$
- (3)  $X_3 = X_2 \cdot 0 = 010, C_3 = 2 \cdot 2 + 0 = 4, L_3 = \lceil \log 3 \rceil + 1 = 3, b_3 = 100$
- (4)  $X_4 = X_0 \cdot 1 = 1, C_4 = 0 \cdot 2 + 1 = 1, L_4 = \lceil \log 4 \rceil + 1 = 3, b_4 = 001$
- (5)  $X_5 = X_1 \cdot 0 = 00, C_5 = 1 \cdot 2 + 0 = 2, L_5 = \lceil \log 5 \rceil + 1 = 4, b_5 = 0010$

따라서, 원 데이터  $X_9 = 0 \cdot 01 \cdot 010 \cdot 1 \cdot 00$ 는  $Y_{13} = 0 \cdot 11 \cdot 100 \cdot 001 \cdot 0010$ 의 압축계열로 변환된다. 위의 예에서 알 수 있듯이 ZL방식의 결점은 입력 데이터의 길이 n이 충분히 길지 않으면 압축의 효과를 기대할 수 없다는 점에 있다. 그 이유는

- (1) 각 부분계열의 마지막 기호  $x_i$ 가 압축되지 않는 형태로 부호화 되고,
- (2)  $i$ 가 커짐에 따라 더 이상 참조되지 않는 부분계열  $X_i$ 가 급격히 증가하기 때문이다. (1)

의 문제점은 부분계열  $X_i$ 로 분해할때 최후의 기호  $x_i$ 를  $X_{i+1}$ 의 선두 기호에 중첩시킴으로써  $j$ 만을 부호화하는 LZW방식에서 해결하고 있다.<sup>(31)</sup> UNIX의 Compress Command<sup>(32)</sup>는 LZW 방식에 근거하고 있다. (2)의 문제점은 2진계열의 경우 증분분해는 완전나무(complete tree)를 구성하여 외부 결점에만 참조번호를 할당시켜 해결할 수 있다.<sup>(20)</sup> 이때 decoding tree에서는 외부 절점의 번호에 의해 데이터 계열을 복호하기 때문에 외부 절점에서 루트 절점의 방향으로 LIFO순서로 데이터가 복원된다. 그러나, Enumerative Code<sup>(33)</sup>의 개념을 도입하여 외부 절점에 번호를 할당하면 FIFO형식으로 동시에 데이터 길이 n에 비례하는  $O(n)$ 의 속도와 소용량으로 복호된다.<sup>(34)</sup>

또한 위의 알고리즘에서는 적응사전을 구성하기 위한 소용 기억 용량에는 제한을 두지 않고 있다. 그래서, 일종의 self-organizing-list를 준비하여 사전의 용량이 포화가 되었을 때 LRU (Least Recently Used) 계열을 제거하게 되면 효율적으로 사용할 수 있다.<sup>(35)</sup> 이외에도 LFU (Least Frequently Used) 계열을 제거하거나, 사전을 Reset, Freeze, Swap하는 등의 방법이 있다.<sup>(36)</sup>

ZL-78방식에서는 Tree나 Hash함수등을 이용하여 부호화하기 때문에 일반적으로 부호화 및 복호화의 처리속도가 고속이나, 증분분해의 부분계열로만 새로운 단어가 사전에 등록되기 때문에 초기에는 압축효율이 좋지 못한 결점이 있다. 한편, ZL-77방식은 과거의 모든 계열을 부호화에 이용할 수 있으나 고속으로 검색, 등록, 소거 등의 처리가 가능한 데이터 구조를 만들기 어렵다. 그래서, 사전의 처리가 고속으로 이루어지도록 Hash함수나 Trie등을 이용하고 동시에 ZL-77방식보다 사전에 등록되는 단어의 수가 증가하도록 고려한 방식이 여러 가지 제안되어 있다.

(37)(38)(39)

V. 산술부호(arithmetic codes)

산술부호는 [0,1]의 數直線 구간을 정보원 기호의 출현확률에 비례하게 순차적으로 분할하여, 정보원 계열(source sequences)을 확률 구간상의 점의 좌표로 부호화하는 방법으로 해석할 수 있다. 산술부호의 기본 개념은 C.E.Shannon<sup>39)</sup>의 논문에서 찾을 수 있으며, Elias의 이상부호(ideal code)<sup>40)</sup>가 제기가 되어 Rissanen과 Pasco에 의해 실용적인 형태로 발전하게 되었다. 그 후 많은 연구가에 의해 개선 발전되어 텍스트 데이터를 비롯하여 영상 데이터의 압축에 이르기까지 다양하게 응용되고 있다.

산술부호는 기존의 잘 알려진 Huffman Codes<sup>38)</sup>와는 달리 블록부호(block code)가 아니며 i.i.d(independently and identically distribution) 정보원에 대해서도 높은 압축율을 달성할 수 있어 장치화의 문제점으로 지적되어 왔던 정보원 확대(alphabet extention)를 수반하지 않는 장점이 있다. 특히, 재귀적 수법(recursive operations)에 의한 알고리즘의 구성이 간단하며, 정보원 모델링과 부호화를 명확히 구분할 수 있어 정보원의 변화에 대한 적응성이 뛰어나므로 다양한 분포의 정보원에도 적용할 수 있는 잇점이 있다.<sup>41)</sup>

5.1 Binary Arithmetic Codes(BAC)

텍스트 중심인 정보원의 캐시밀러 통신에는 MH(Modified Huffman), MR(Modified READ) 및 MMR방식이 주로 사용되어 왔으나, 최근의 문서통신의 다양화에 따라 등장된 몇가지 방식중에 산술부호를 이용한 방식이 주목되고 있다. 여기서는 BAC 알고리즘을 중심으로 소개하고<sup>42)</sup> 그 응용에 대하여 알아 본다.

길이 n의 정보원 계열  $X_n = x_1x_2 \dots x_n$ ,  $x_i \in \{0,1\}$ 의 부호화를 생각한다. 기호 0,1의 출현확률이  $p(0)$ ,  $p(1)$ 일때, MPS(More Probability Symbol)를 0, LPS(Less Probability Symbol)를 1로 가정한다. 산술부호의 확률구간[0,1]을 아래와 같이 q 비트의 2진 소수로 나타낸다.

$$A(\lambda) = \underbrace{0.111 \dots 1}_{q \text{ bits}} \quad (1)$$

길이 n-1의 계열  $X_{n-1} = x_1x_2 \dots x_{n-1}$ 의 부호화가 끝난 상태에서  $X_n$ 의 부호화를 위한 부분구간은 (2)식에 의해 분할된다.

$$\begin{aligned} A(X_{n-1} \cdot 1) &= A(X_{n-1}) \times p(1) \\ A(X_{n-1} \cdot 0) &= 1 - A(X_{n-1} \cdot 1) \end{aligned} \quad (2)$$

이때 계열  $X_n$ 의 부호어  $C(X_n)$ 은 (3)식에 의해 구해진다.

$$C(X_n) = \begin{cases} C(X_{n-1} \cdot 0) = C(X_{n-1}), & x_n = 0 \\ C(X_{n-1} \cdot 1) = C(X_{n-1}) + A(X_{n-1} \cdot 0), & x_n = 1 \end{cases} \quad (3)$$

한편, 계열  $X'_n$ 의 복호가 완료된 상태에서  $X'_n = X'_{n-1} \cdot x'_n$ 는  $x'_n$ 는

$$\begin{aligned} C(X'_n) < C(X'_{n-1}) < A(X'_{n-1} \cdot 0) \text{이면, } x'_n &\rightarrow 0 \\ \text{아니면, } x'_n &\rightarrow 1 \end{aligned} \quad (4)$$

로 복호화가 이루어진다.

이상에서 본 바와 같이 산술부호의 계산량은 구간 분할에 필요한 곱셈에 지배되며, 순차적으로 구간폭 A가 좁아지게 되므로 이진소수의 精度(precision)가 증가하는 문제가 발생한다. 이것을 해결하기 위해 Rissanen은 LPS의 확률을  $2^{-Q}$  ( $Q : \text{integer}$ )로 근사화시켜 곱셈을 쉬프트 조작으로 대신하는 방법을 제안하고 있다.<sup>43)</sup>

$$p(1) \approx p(1)' = 2^{-Q}, \quad 1 : \text{LPS}, \quad Q : \text{skew number}$$

$$\begin{cases} A(\lambda) = 0.111 \dots 1, & q \text{ bits precision} \\ C(\lambda) = 0.000 \dots 0 \end{cases} \quad (5)$$

$$A(X_{n-1} \cdot 1) = \langle A(X_{n-1}) \cdot 2^{-Q} \rangle_{q-Q} \quad (6)$$

$$A(X_{n-1} \cdot 0) = A(X_{n-1}) - A(X_{n-1} \cdot 1) \quad (7)$$

여기서,  $\langle X \rangle_q$ 는 2진수 X의 유효 비트수 q비트로 절삭함을 의미한다. 길이 8인 계열  $X_8 = 01001$



101, (Q=2, q=4, p(1)=1/4, p(0)=3/4)의 산술부호화의 과정을 표2에 보인다.

표 2. BAC에 의한 산술 부호화의 예  
Table 2. Arithmetic coding example using BAC.

$x_n$	$A(X_{n-1})$	$A(X_{n-1}, 0)$	$A(X_{n-1}, 1)$	$C(X_{n-1}, x_n)$
0	.1111	.1100	.0011	.0000
1	.1100	.1001	.0011	.1001
0	.001100	.001001	.000011	.1001
0	.001001	.000111	.000010	.1001
1	.0001110	.0001011	.0000011	.1010011
1	.000001100	.000001001	.000000011	.101010101
0	.00000001100	.00000001001	.00000000011	.101010101
1	.000000001001	.000000000111	.000000000010	.10101011011

산술부호의 효율(efficiency)은 최종 부분구간의 좌단점  $C(X_n)$ 을 나타내기 위해 필요한 비트 수로 주어지는 평균부호 길이(average code length)로 평가된다.  $C(X_n)$ 의 유효 자릿수는  $A(X_n)$ 의 자릿수와 동일하며, 정보원 계열  $X_n$ 에서 LPS가 출현하면 Q만큼 증가하게 된다.  $(1-2^{-Q})^d \leq 1/2$ 을 만족하는 d에 대해,  $X_n$ 중에 MPS가 d번 나타나면 유효 자릿수는 한 자릿수 증가한다. 길이가 충분히 큰 n에 대하여 LPS는 np, MPS는  $n(1-p)$ 회 출현하므로  $A(X_n)$ 의 소숫점 이하의 자릿수는,

$$npQ + \frac{n(1-p)}{d} = npQ - n(1-p)\log(1-2^{-Q}) \quad (8)$$

로 된다. 따라서 한 기호당 평균부호길이  $L_n = pQ - (1-p)\log(1-2^{-Q})$ 로 되어, 유효 자릿수 q를 크게 취하고  $n \rightarrow \infty$ 일때,  $L_n \rightarrow H(X_n)$ 에 접근하는 점근적 최량성(asymptotical optimality)이 보장되는 부호이다. 여기서,  $H(X_n)$ 은 계열  $X_n$ 의 엔트로피이다.

BAC를 이용한 B/W FAX 영상의 압축효과는 G3 FAX의 표준방식인 MR(Modified READ)에 비해 20~30%의 압축율이 향상됨을 보이고 있어 많은 응용이 전개되고 있다.<sup>(43)</sup> 특히, CCITT와 ISO를 중심으로 한 JPEG(Joint

Picture Expert Group) 및 JBIG(Joint Bi-level Image Group)의 표준화 활동에 의해 크게 발전하게 되었다. JPEG와 JBIG의 엔트로피 부호화(entropy coding)을 위해 IBM 중심의 Q-coder와 MITSUBISHI의 MEL coder 그리고 AT&T의 Minimax coder의 탄생을 보게 되었다.<sup>(44)</sup> 이들은 모두 산술부호의 기본 개념을 기초로 하여 실용화를 위해 설계된 부호들이다. Q-coder<sup>(45)</sup>는 기호의 출현확률을 평가하여 확률 적응화를 실행하고 있다. Q coder에서는 LPS의 출현확률인 Q의 대표점이 미리 준비되어 있어서 이 Q를 가리키는 인덱스를 적응적으로 움직임으로써 동적 부호화(dynamic coding)가 가능하도록 한 특징을 갖고 있다.

### 5.2 Multi-level Arithmetic Codes(MAC)

FAX의 적용분야가 확대됨에 따라 종래의 문자나 그림과 같은 2진영상(binary image)만을 취급해온 FAX도 사진(photographic)과 같은 멀티 레벨 및 칼라의 영상 데이터를 취급하기에 이르렀다. 이와 같은 FAX신호의 부호화에는 자연 화상의 부호화와는 다른 무잡음형 부호화(noiseless coding)가 주류를 이루고 있으며, 정보원의 국소적인 성질(locally statistics)을 잘 반영하는 산술부호의 이용이 좋은 효과를 얻을 수 있다.

멀티 레벨 산술 부호는 BAC의 확장으로 쉽게 실현될 수 있다. 정보원 계열  $X_n = x_1 x_2 \dots x_n$ ,  $x_i \in A = \{0, 1, \dots, k-1\}$ ,  $p(0) \geq p(1) \geq \dots \geq p(k-1)$ 에 대한 부호화는 아래와 같이 이루어진다.

$$A(\lambda) = 0.11 \dots 1, \quad q \text{ bits precision} \quad (9)$$

$$A(x_i) = A(\lambda) \cdot p(x_i), \quad x_i \in A \quad (10)$$

$$C(x_i) = \begin{cases} \sum_{l=1}^{x_i} A(l), & x_i \neq 0 \text{ (1st symbol cordword)} \\ 0, & x_i = 0 \end{cases} \quad (11)$$

$$A(X_n) = A(X_{n-1}, x_n) = A(X_{n-1}, 1) \cdot p(x_n), \quad (12)$$

$$C(X_n) = \begin{cases} C(X_{n-1}) + \sum_{l=1}^k A(X_{n-1} \cdot l), & x_n \neq 0 \\ C(X_{n-1}), & x_n = 0 \end{cases} \quad (13)$$

여기서 (10), (12)식의 곱셈은 다음과 같이 쉬프트 조작으로 대체하여 연산을 고속화시킬 수 있다.

$$A(l) = \langle A(\lambda) \cdot 2^{ql} \rangle_{q, q_1}, \quad 0 \leq l \leq k-1 \quad (14)$$

$$A(0) = A(\lambda) - \sum_{l=1}^{k-1} A(l) \quad (15)$$

$$A(X_{n-1} \cdot 1) = \langle A(X_{n-1}) \cdot 2^{q_1} \rangle_{q, q_1} \quad (16)$$

$$A(X_{n-1} \cdot 0) = A(X_{n-1}) - \sum_{l=1}^{k-1} A(X_{n-1} \cdot l) \quad (17)$$

따라서, 멀티 레벨 산술부호의 평균부호 길이 L은 (18)식으로 주어진다.

$$L = \sum_{l=1}^k \{p(l) \times Q_l\} - p(0) \log(1 - \sum_{l=1}^k 2^{-q_l}) \quad (18)$$

산술부호를 실제로 구현함에 있어서의 문제점은 精度의 증가 및 자리올림(carry over)의 문제가 일어나므로 이를 위한 해결법이 상구되어야 한다. 精度증가의 문제는 쉬프트 조작(shift operation)에 의해 간단히 해결되며, 자리올림 문제는 그림3의 비트 삽입법(bit stuffing)<sup>(43)</sup>으로 해결할 수 있다. 자리올림은  $C(X_{n-1})$  및  $A(X_{n-1} \cdot x_n)$ 을 q개의 비트로 표시하여 덧셈 조작을 할 때 일어난다.  $C(X_{n-1})$ 과  $A(X_{n-1} \cdot x_n)$ 의 덧셈으로 얻어지는 부호 계열(code string)은 상위 레지스터 V 및 M을 거쳐 전송모로 보내지게 되나, 그림3과 같이 V레지스터의 비트가 전부 1로 채워져 있다면 그 상위의 M레지스터까지 자리올림의 영향을 주게되어 많은 비트 반전을 요구하게 된다. 더욱 파급되는 자리 올림은 전송이 끝난 부호 계열에까지 영향을 미치게 되므로 정상적인 정보전달이 불가능 하게 된다. 이를 위한 해결 방안의 하나로 Rissanen은 비트 삽입법을 제안하였고, MORITA<sup>(46)</sup>는 레지스터 V와 M에 걸쳐서 생기는 자리올림을 해결하기 위하여

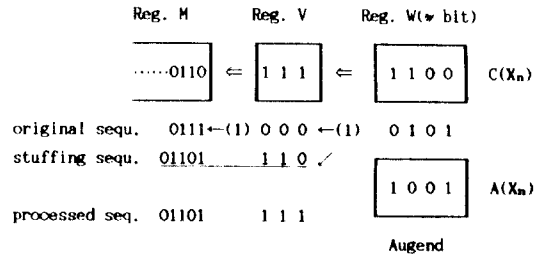


그림 3. 비트 삽입법의 일례  
Fig 3. bit stuffing example.

비트 삽입법을 확장하는 방법을 제시하였다.

### 5.3 MAC의 고속화

데이터 압축의 또다른 목표는 알고리즘의 고속화이다. 높은 압축효율을 얻을 수 있는 산술부호의 계산량은 입력 데이터의 길이 n에 선형적으로 비례하는  $O(n)$ 으로 개선되어 실용적으로 응용 가능하게 되었다. 그러나, 많은 산술 연산 때문에 적응 사전법의 유니버설 압축법에 비하여 처리속도가 대단히 느린 결점이 있다.<sup>(44)</sup> 한편, 산술부호와 같은 모델 베이스 유니버설 압축법에서는 적응적으로 화물 파라메타를 생성해야 하기 때문에 모델링과 부호화를 함께 고려해야 바람직하다.

자분모델<sup>(47)</sup>과 산술 부호의 일종인 Pasco의 FV부호(Fixed to Variable codes)<sup>(48)</sup>를 이용하여 gray level image를 압축시의 산술부호화를 고속화하는 한 방법을 보인다.<sup>(49)</sup> Pasco의 산술부호는 Elias부호의  $O(n^2)$ 의 계산량을  $O(n)$ 으로 개선시킬 것으로 알고리즘이 간단하며 Elias부호를 가장 충실히 실현하는 부호이다. 길이 n의 정보원 계열  $X_n = x_1 x_2 \dots x_n$ 에 대하여 각 기호  $x_i \in A = \{0, 1, \dots, k-1\}$ 의 출현확률  $p(x_i)$ 를 J비트 이하의 2진 小數로 나타낸다. 이때  $p(x_i)$ 는 다음의 조건을 만족해야 한다.

$$0 \leq p(x_i) < 1, \quad \sum p(x_i) = 1 \quad (19)$$

또한 각 기호  $x_i$ 에 대한 누적확률  $C(x_i)$ 는

$$c(\lambda)=0 \quad (20)$$

$$c(x_i)=\sum_{j=1}^n p(x_j) \quad (21)$$

$$c(x_i)+p(x_i)\leq 1 \quad (22)$$

로 된다. Pasco의 산술부호화 알고리즘은 다음과 같다.

<encoding process>

- (1)  $F \rightarrow 0$
- (2)  $T \rightarrow 1$
- (3) For  $i=1$  to  $n$  do
  - (3.1) Load the  $p(x_i)$  and  $c(x_i)$
  - (3.2)  $F_i \rightarrow F_{i-1} + T_{i-1} \times c(x_i)$
  - (3.3)  $T_i \rightarrow T_{i-1} \times p(x_i)$ , (truncated to  $q$  significant digits)
- (4)  $L \rightarrow \tau + 1$  ( $T_n = t_0 t_1 t_2 \dots t_{q-1} \times b^\tau$ ,  $t_0 \neq 0$ ), ( $\lceil -\log_b T_n \rceil + 1$ )
- (5)  $F_n \rightarrow F_n + b^L$
- (6) transmit the  $L$  digits of  $F_n$ . (code string

$$Y_L = y_1 y_2 \dots y_L$$

위 알고리즘에서 각 변수( $F, T, c, p$ )는 小數이지만 산술연산을 고속으로 하기 위하여 모두 정수를 이용한다. 이를 위해 출현확률  $p(x_i)$ 와 누적확률  $c(x_i)$ 를  $2^{16}$ 배 하므로 (19)식의 총합은  $2^{16}$ 이 된다. 또  $T$ 와  $F$ 의 값은 부호화와 동시에 점점 작아지나  $T$ 의 유효 자릿수를  $q=16$ 으로 하여  $F$ 와 함께 순차적으로 shift out한다. 그러나  $T$ 와  $F$ 는 알고리즘(3.2)에서 곱셈과 덧셈이 이루어지므로 32비트의 정수를 이용한다.  $F$ 로부터 shift out되는 상위의 비트가 해당 부호어으로써 출력되지만 자리 올림 문제를 해결하기 위하여 비트 삽입법을 부가하고 있다.

또한, 멀티 레벨 산술 부호의 실제의 연산속도는 encoding process의 (3)에 크게 의존한다. (3.1)의  $p(x_i)$ 와  $c(x_i)$ 는 정보원 모델링의 갱신과 함께 동적으로 계산되는 량이므로 거의 일정한 파라메타이나,  $F$ 와  $T$ 는

$$F_i = \begin{cases} F_{i-1} + \sum_{j=1}^{y-1} T_{i-1} \times p(j), & y \neq 0 \\ F_{i-1}, & y = 0 \end{cases} \quad (23)$$

$$T_i = T_{i-1} \times p(y), \quad y : x_i \text{의 gray-level value} \quad (24)$$

로 구해진다. 적응 산술부호화에 필요한 연산량은,

- (1) 출현확률  $p(x_i)$ 를 위한 나눗셈
- (2) 부분구간  $F_i$ 를 위한 곱셈 및 덧셈
- (3) 각 부분구간을 누적하기 위한 덧셈 및 뺄셈

등이다. 따라서 상대적인 연산량을 줄이기 위해서는 높은 확률의 기호에 해당하는 부분구간일수록 분할의 시점 가까이에 위치 시키면 된다. 이 착상을 근접 화소의 차분값이  $y=0$  부근에 집중하는 차분모델에 적용시키는 것은 용이하다. 즉 차분값을 확률구간의 시작점으로부터  $0, -1, 1, \dots, -(k-1), (k-1)$ 의 순서로 설정하면 필요한 연산량이 대폭 줄어들게 되므로 고속처리가 가능하게 된다.

## VI. 결 론

이 논문에서는 다양한 정보원에 적용할 수 있는 유니버설 부호의 기본 개념을 보이고 그 응용을 중심으로 한 최근의 연구 동향에 대하여 알아보았다. 문서통신의 표준방식으로 사용되어 온 MH부호의 기본인 Huffman부호를 적응적으로 처리하는 새로운 방식의 등장과 자연수를 효과적으로 표현하기 위한 유니버설부호는 데이터 압축시스템의 통합적인 구축을 위해 효과적인 방법이 될 것이다.

적용 사전법의 대표적인 방식인 Ziv-Lempel 부호를 알고리즘 중심으로 나타내었으며 각종 개량방식 및 실용화 방식에 대하여 알아보았다. 현재 PC통신의 압축 유틸리티의 주류를 이루고 있는 PKARC, LZH, ZIP 및 LHARC는 Ziv Lempel부호의 응용인 점을 주시하면 보다 활발한 연구 환경의 조성으로 우리의 고유

한 알고리즘의 개발이 시급한 실정이라 하였다. 또한 FAX를 비롯한 영상 데이터의 압축을 위한 모델과 산술부호에 대하여 그 기본 원리로부터 확장 및 응용에 이르기까지 알아보았다. 그러나 여기서 소개한 많은 부호들 가운데 그 성능에 관하여 이론적으로 또는 시뮬레이션등으로 자세하게 해석 평가된 것은 극소수에 지나지 않는다. 따라서 각 부호의 개량법등에 비교적 쉽게 얻어질 수도 있으므로 이 분야의 중요자료로 활용될 것을 기대한다.

### 참 고 문 헌

1. T. Berger., "Rate distortion theory", prentice Hall(1971).
2. C.E.Shannon., "the mathematical theory of communication", University of Illinois Press, Urbana (1949).
3. D.A.Huffman., "A method for the construction of minimum-redundancy codes", Proc. IRE, 40, 9, pp. 1098-1101 (Sept. 1952).
4. R.E.Blahut., "Principles and practice of information theory", Addison-wesley.
5. D.A.Lewler and D.S.Hirschberg., "Data Compression", ACM Computing surveys, 19, 3, pp.261-296 (Sep. 1987).
6. J.Rissanen., "Generalized Kraft inequality and arithmetic coding", IBM J.Res.Dev.20,pp.198-203 (May 1979).
7. D.E.Knuth., "Dynamic Huffman coding", J.Algorithms, 6,2,pp.163-180 (Jun. 1985).
8. A.Lempel, J.Ziv., "On the complexity of finite sequences", IEEE Trans. IT 22, 1, pp.75-81 (Jan. 1976).
9. L.Bentley et al., "A locally adaptive data compression scheme", CACM 29, 4, pp.320-330(Apr.1986).
10. F.M.Willems., "Universal data compression and repetition times", IEEE Trans. IT 35, 1,pp.54-58 (Jan. 1989).
11. J.H.PARK and H.IMAL., "Data compression", KITE 15, 1,pp.99-105 (2 1988).
12. N.Faller., "An adaptive system for data compression", Record of the 7 th Asilomar Conf. on Circuits, System and Computers, pp.593-597, (1973).
13. R.G.Gallager, "Variations on a theme by Huffman", IEEE Trans. IT 24, 6,pp.668-675 (Nov.1978).
14. J.S.Vitter., "Design and analysis of dynamic Huffman coding", JACM 34, 4,pp.825-845 (Oct.1987).
15. E.S.Schwartz., "An optimum encoding with minimum longest code and total number digits", Inf. Contorl , 7,1,pp.37-44 (mar. 1964).
16. C.L.McMaster., "Documentation of the compact command", In UNIX User's Manual, 4.2 BSD, Virtual VAX 11 Version, Univ.of Califo, Berkeley (Mar. 1981).
17. 김철배, 민용식, "효율적인 한글코드화에 관한 연구", 한국통신학회논문지 14, 6, pp.633-641 (12.1989)
18. 김현진 외 3인, "히프만 부호를 이용한 압축/인출 결합 알고리즘", 한국정보과학회 '91 봄학술발표논문집 18.1 (4.1991).
19. J.A.Storer., "Data compression : methods and theory", Compu. Sci. Press(1988).
20. H.Yamamoto et al., "Ziv Lempel符號の改良とツミエシツェンによる性能評價", IEICE Tech. Report CS84-135, pp.1-8 (1.1985).
21. P.Elias., "Universal codeword sets and representations of the integers", IEEE Trans IT-21, 2, pp.194-203 (Mar.1975).
22. J.H.PARK et al., "An adaptive data compression using self organizing rule", IEICE Trans J72-A, 8 ,pp.1353-1359 (8.1989).
23. A.Apostolico et al., "Robust transmission of unbounded strings using Fibonacci representations", IEEE Trans IT 33, 2, pp.238-245 (Mar.1987).
24. H.Yamamoto et al., "有限正整数に対する2値符號化法の検討", SITA 12, 2, pp.559-564 (12.1989).
25. J.Ziv and A.Lempel., "An universal algorithm for sequential data compression", IEEE Trans IT 23, 3, pp.337-343 (May. 1977).
26. J.Ziv and A.Lempel., "Compression of individual sequences via variable rate coding", IEEE Trans. IT 24, 5, pp.530-536 (Sept.1978).
27. PKARC FAST., "PKARC FAST File Archival Utility", Version 3.5 PKWARE, Inc. Glendale, Wis. (1987).
28. D.E.Knuth et al., "Fast pattern matching in strings", SIAM. J. of comput., Vol.6, No.2 pp.323-350 (June, 1977).

29. M.Rodeh et al., "Linear algorithm for data compression via string matching", J. of ACM, 28, 1, pp. 16-24 (Jan.1981).

30. J.H.PARK et al., "Practical Data Compression Methods using Pattern Matching and Arithmetic Codes", IEICE Trans J71-A, 8, pp.1615-1623 (1988, 8, 8).

31. T.A.Welch., "A technique for high performance data compression", IEEE computer, pp.8-19 (June.1984).

32. S.W.Thomas et al., "Compress(version 4.0) program and documentation", available from joepetsd, UUCP.

33. T.M Cover, "Enumerative source coding", IEEE Trans IT-19.1, pp.73-77 (1.1973).

34. T.Kawabata et al., "A new implementation of Ziv Lempel incremental parsing algorithm", submitted to IEEE Trans Information Theory.

35. R.N.Horspool., "A locally adaptive data compression scheme", CACM, 30,9, pp.792 (Sept.1987).

36. J.A.Storer., "Parallel algorithms for On line data compression", IEEE ICC'88, pp.385-389 (1988).

37. E.R.Fiala et al., "Data Compression with finite windows", CACM, 32, 4, pp.490-505 (April.1989).

38. 村島,中村,廣田., "文字列登録を強化したLZW法によるソフトウェアの圧縮", IEICE Trans J73-A, 9, pp. 1529-1533 (Sept.1990).

39. 森田,小林., "再生可能な文字列分解にもとづくユニバーサルデータ圧縮について", SITA, 13, 2(1.1991).

40. N.Abramson., "Information theory and coding", New York McGraw-Hill, 1963.

41. J.Rissanen and G.G.Langdon., "Universal modeling and coding", IEEE Trans. IT-27, 1, pp.12-23 (Jan. 1981).

42. G.G.Langdon., "An introduction to arithmetic coding", IBM J.Res. Dev., 28, 2, (1984).

43. G.G.Langdon and J.Rissanen., "Compression of B/W images with arithmetic coding", IEEE Trans. COM-29, 6, pp.858-867 (1981).

44. H.Yasuda., "International sandardization of high efficiency coding", J. ITE, 43, 10, pp.1011-1019, (1989).

45. Q-Coder Special, "Q-Coder adaptive binary arithmetic coder", IBM J. Res. Dev. 32, 6, pp.715-840 (Nov. 1989).

46. H.Morita et al., "On efficiency of binary arithmetic codes", J. IPA, 25, 4, pp.622-631 (1984).

47. J.H.Park and H. Imai., "Arithmetic coding for gray-level image data compression", PCSJ 89, 7-18, (1989).

48. Pasco R.C., "Source coding algorithm for fast data compression", Ph.D Dissertation Stanford Univ. (1976).



林志煥(Ji Hwan PARK) 正會員  
 1958年 2月 6日生  
 1987年 3月: 日本電氣通信大學 應用電子工學科(工學修士)  
 1990年 3月: 日本立山大學國立大學 電子情報工學科(工學博士)  
 1990年 3月~現在: 釜山水産大學校 電子計算學科 專任講師



陳 庸 玉(Yong Ohk CHIN) 正會員  
 1943年 3月 21日生  
 1968年 2月: 延世大學校 電氣工學科 卒業(工學士)  
 1975年 2月: 延世大學校 大學院 電子工學科(工學碩士)  
 1981年 2月: 延世大學校 大學院 電子工學科(工學博士)  
 1975年~1978年: 光云工科大学 通信工學科 教授  
 1980年~現在: 慶熙大學校 電子工學科 教授  
 1980年 通信技術士