

グラフマッチングによる
文書画像のストローク抽出法
道庭賢一 佐藤 誠
東京工業大学精密工学研究所

あらまし 情報化社会の高度化に伴い、大量の文書画像データを効率よく処理することがますます重要な問題となっている。伝達、理解、管理など文書画像の処理は従来別々に扱われてきたが、これらは本来統合して処理されるべきものである。このためには文書画像を階層的に表現することが必要であり、ここでは階層表現に適した新しい2値画像の表現法、最小被覆ラン表現(MCR表現)を提案する。これは2値画像を最小個数の縦ランと横ランで表現する方法である。本研究では、MCR表現を求める問題が2部グラフの最大マッチング問題に帰着することを示すとともに、効率のよいアルゴリズムを明らかにした。

A S t r o k e E x t r a c t i o n M e t h o d
o f D o c u m e n t I m a g e s
B a s e d o n M a t c h i n g o f B i p a r t i t e G r a p h

Kenichi Douniwa Makoto Sato

Research Laboratory of Precision Machinery and Electronics
TOKYO INSTITUTE OF TECHNOLOGY

Abstract With the advance of information-oriented society, it is necessary to treat a great deal of document images efficiently. Document image processing, such as coding, understanding, and managing, have been treated separately. But the integration of document image processing, based on hierarchical structures of document images, is desirable. In this paper we propose a new expression of binary images, called Minimum Covering Run Expression. The expression is suitable for hierarchical description of binary images. It is shown that the expression method is closely related to the maximum matching problem of bipartite graph. The efficient algorithm for the expression method is also demonstrated.

1 まえがき

情報化社会の発展にともなって膨大な量の文書画像が取り扱われるようになり、この文書画像データを効率よく処理することが重要な課題となっている。

文書画像の基本的な処理としては、ファクシミリなどによる伝達、文字領域・図形領域・画像領域への領域分割や文字認識等の文書理解、データの蓄積や検索によるデータベース管理などが考えられる。これらの伝達、理解、管理などの処理は従来独立して処理されてきたが、本来統合して処理されるべきものである。例えば、CCITT（国際電信電話諮問委員会）においてすでに標準化されているMH法、MR法は、ファクシミリにおける伝達のための効率の良い符号化であるが、理解、管理に関しては必ずしも十分に考慮されていない。

文書画像の構造は図1のように“ページ”、“フレーム”、“ブロック”、“ストローク”、“ラン”、“画素”という構成要素を用いて階層的に表現できる。文字、図形、イメージからなる文書画像を、このような階層構造により表現するための標準化の検討がCCITTとISOにおいて行なわれており、ODA (Open/Office Document Architecture)⁽¹⁾⁽²⁾と呼ばれている。また階層構造を利用した文書処理システムの検討も盛んに行われている⁽³⁾⁽⁴⁾。この階層構造に基づいた表現が実現すれば、伝達、管理のための符号化と、理解のための解析を統合して処理することができる。

文書画像の処理の統合を考えた場合、入力される文書画像から自動的にその階層的表現を得る問題が重要である。2値画像の基本的な表現法としては、ビットマップ表現、横ランを要素とした表現などがある。しかし、これらは階層構造を表現することは考慮されていない。ここでは階層表現に適した新しい2値画像の表現法として、最小被覆ラン表現 (Minimum Covering Run Expression) を提案する。以下ではこれをMCR表現とも呼ぶ。MCR表現は縦ランと横ランを基本要素として、2値画像を最小個数のランで被覆することにより表現する方法である。2値画像の縦ランと横ランの集合に、2部グラフ (Bipartite Graph) の構造を導入することにより、MCR表現を求める問題は、2部グラフ

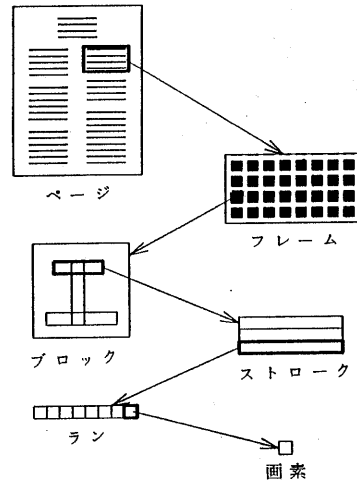


図1 文書画像の階層表現

の最大マッチング問題に帰着できる。

MCR表現の縦ラン被覆部分、横ラン被覆部分のセグメントは、それぞれ縦ストローク、横ストロークに良く対応するため、階層的表現に適し、文書画像の構造解析にも有効である。

2章で、2部グラフとその最大マッチングに関する準備を行なう。3章で2値画像のMCR表現を定義する。4章でMCR表現のための解析アルゴリズムを明らかにし、5章で実験結果を示す。

2 2部グラフと最大マッチング

2.1 2部グラフ

頂点の集合を $V(G)$ 、辺の集合を $E(G)$ とする。互いに素な集合 X と Y に対して

$$V(G) = X \cup Y$$

$$E(G) = \{ \langle x_i, y_j \rangle \mid x_i \in X, y_j \in Y \}$$

となるグラフを X と Y を部集合とする2部グラフという。ただし $\langle x_i, y_j \rangle$ は頂点 x_i, y_j の間の辺を表す。2部グラフは隣接行列により表現することができる。すなわち2部グラフ G の隣接行列

$$F(X, Y) = [f_{ij}]$$

を次のように定義する。 $F(X, Y)$ は各要素 f_{ij} を

$$f_{ij} = \begin{cases} 1: \langle x_i, y_j \rangle \in E(G) \\ 0: \text{そうでないもの} \end{cases}$$

とする $|X| \times |Y|$ 行列である。ただし

$|X|$, $|Y|$ はそれぞれ集合 X , Y の要素の数である。2部グラフとその隣接行列の例を図2(a), (b)に示す。

2. 2 最大マッチングと最小点被覆

辺の集合 $M \subseteq E$ において、どの相異なる2つの辺 $e_i, e_j \in M$ も端点を共有しないとき、 M を2部グラフ G のマッチングという。つまり行列 $F(X, Y)$ の各行、各列において非零要素をたかだか1個選んだとき、この集合 M がこの2部グラフのマッチングとなる。2部グラフ G のマッチングのうちで、辺の数が最大のものを G の最大マッチングという。図2の2部グラフの最大マッチングは、図2(a)の太い線に対応する辺の集合である。図2(b)では丸で囲まれた非零要素の集合がそれに対応する。したがって、この場合の最大マッチング M は

$$M = \{ \langle x_1, y_2 \rangle, \langle x_2, y_1 \rangle, \langle x_3, y_3 \rangle \}$$

となる。

次に頂点の集合 $W \subseteq V$ がどの辺 $e \in E$ の端点をも少なくとも1個含むとき、 W を2部グラフ G の点被覆という。つまり行列 $F(X, Y)$ のすべての非零要素を被覆するように行および列に直線を引いたとき、この直線に対応する頂点が G の点被覆になる。そして2部グラフ G の点被覆のうちで頂点の数が最小のものを最小点被覆という。図2の2部グラフの最小点被覆は、図2(a)の黒点に対応する頂点の集合である。図2(b)では、行および列に引いた直線に対応する頂点の集合がそれに対応する。したがって最小点被覆 W は

$$W = \{ x_2, x_3, y_2 \}$$

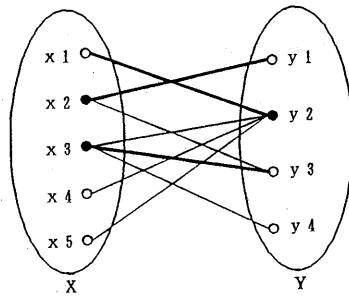
となる。

2部グラフの最大マッチング M , 最小点被覆 W に対して、 $|M| = |W|$ が成り立つことが König の定理より導かれる⁽⁵⁾。そして2部グラフの最大マッチングを求める問題と最小点被覆を求める問題は互いに等価である。

2部グラフの最大マッチング問題は、一般グラフに対する場合よりもはるかに効率的に求めることができる。例えば Hopcroft-Karp は

$$O(\sqrt{n} \cdot (m + n))$$

ただし $m = |E|$, $n = |V|$



(a)

$$F(X, Y) = \begin{matrix} & y_1 & y_2 & y_3 & y_4 \\ \begin{matrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \end{matrix} & \begin{pmatrix} \cdot & \cdot & \cdot & \cdot \\ \textcircled{1} & \cdot & 1 & \cdot \\ \cdot & 1 & \textcircled{1} & 1 \\ \cdot & 1 & \cdot & \cdot \\ \cdot & 1 & \cdot & \cdot \end{pmatrix} \end{matrix}$$

(b)

図2 2部グラフと隣接行列

の手続き数で求めるアルゴリズムを提案している⁽⁶⁾。

3 2値画像の最小被覆ラン表現

3. 1 最小被覆ラン表現 (MCR表現)

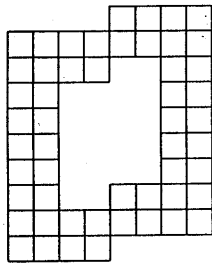
2値画像の横ランを集合 X の要素とし、縦ランを集合 Y の要素とする。互いに交差し合うランの対に2部グラフの辺を対応させると、縦ランどうし横ランどうしは交差ししないため、辺は必ず一方の端点を X の集合内に、他方の端点を Y の集合内に持つ。このようにして2値画像に2部グラフを導入することができる。

MCR表現とは、縦ランと横ランを基本要素として、最小個数のランで黒画素領域を被覆することにより2値画像を表現する方法である。したがってMCR表現を求める問題は、2部グラフの最小点被覆問題、あるいは最大マッチング問題に帰着することができる。

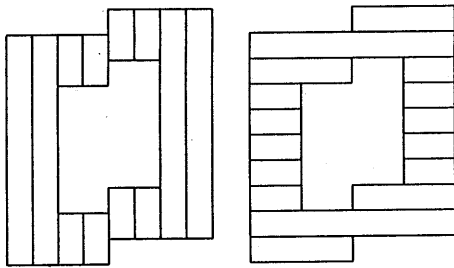
3. 2 MCR表現とストローク要素

文書画像の多くはストロークを基本要素として構成されている。文書画像中の縦ストロークは横ランで被覆するより、縦ランで被覆するほうがランの数が少なくなる。横ストロークも同様に、横ランで被覆するほうがランの数が少なくなる。

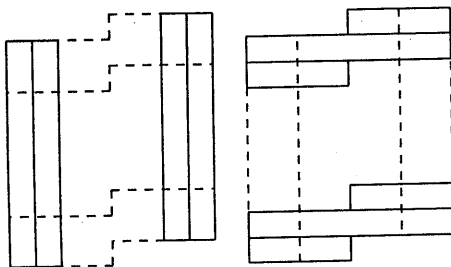
逆に、MCR表現の縦ラン被覆部分、横ラン被覆部分のそれぞれの連結要素であるセグメントは、縦ストローク、横ストロークに良く対応していると考えられる。



(a) 2値画像パターン



縦ラン要素 横ラン要素
(b) 2値画像パターンのラン要素



縦ラン被覆部 横ラン被覆部
(c) MCR表現

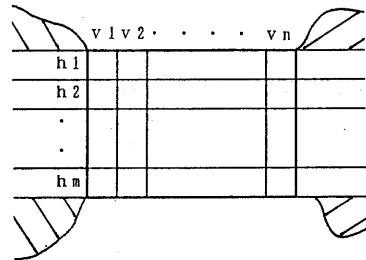
図 3

具体例を図3(a)に示す。このパターンの縦ランと横ラン要素を図3(b)に示す。このMCR表現を求め、縦ラン被覆部と横ラン被覆部をそれぞれ図3(c)に示す。縦ラン被覆部、横ラン被覆部の各セグメントがそれぞれ縦ストローク、横ストロークに対応していることが分かる。

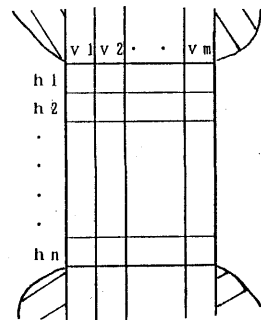
3. 3 MCR表現のための矩形解析

文書中に大きな表やグラフなどがある場合、これらの部分は黒画素の連結成分の数が一般には大きなものとなる。この部分のMCR表現を得るには、大規模な2部グラフの最大マッチングを求めることになり、処理時間が問題となる。そのため、ランデータを登録する前処理の段階で、局所的に形状を判断することにより、予めMCR表現の要素になるラン(被覆ラン)と、ならないラン(非被覆ラン)をある程度確定できると良い。

図4(a), (b)のような矩形を含む領域を考えることにする。このとき次の定理が成り立つ。



(a)



(b)

図 4 最大マッチングの矩形解析

[定理] 図4(a)のように上辺と下辺を境界とする縦 m 、横 n ($m \leq n$)の横長の矩形を含む領域について、横ラン h_1, h_2, \dots, h_m はMCR表現の被覆ランで、縦ラン v_1, v_2, \dots, v_n は非被覆ランである。

同様に図4(b)のように、左辺と右辺を境界とする縦 n 、横 m ($m \leq n$)の縦長の矩形を含む領域について、縦ラン v_1, v_2, \dots, v_m はMCR表現の被覆ランで、横ラン h_1, h_2, \dots, h_n は非被覆ランである。

この定理により図4(a), (b)のような矩形領域を調べることににより、前処理の段階で多くのランを被覆ラン、非被覆ランに確定できる。そして最大マッチング問題の対象となる領域は、図4の斜線部で示した部分領域である。これらの領域はもとの2値画像に比べてはるかに細分化されているので、処理の高速化が期待できる。この矩形領域の局所的処理を、以後MCR表現のための矩形解析と呼ぶ。

4 MCR表現の解析アルゴリズム

4.1 ランデータ構造

最大マッチングで必要になる隣接行列を生成するためには、交差するランが効率よく求められるように縦ラン、横ランのデータを管理する必要がある。横ランデータの管理方法は図5(b)のようにする。

横ランデータを管理するために、二つの配列を用意する。配列 $h_{ran}()$ は、横ランデータの始点座標(sp)、終点座標(ep)、フラグ(fg)、および同一走査線状の次のランデータの保存されている番地を表すポイント(np)から構成される。フラグ(fg)は、被覆ラン、非被覆ラン、あるいはいずれにも確定していない未処理ランなどのランの属性を表現する。 $h_{top}()$ は、各水平走査線に属する最初の横ランデータのポイントを設定する。この横ランのデータ構造により、任意の水平走査線上のランデータを効率よく取り出すことができる。

縦ランのデータ構造も、全く同様にして2つの配列 $v_{ran}()$ と $v_{top}()$ により表現される。

このようなランデータ構造を用いることによ

り、交差するランは次のようにして求めることができる。

垂直走査線 y 上で、始点座標が x_s 、終点座標が x_e の縦ランと交差する横ランは

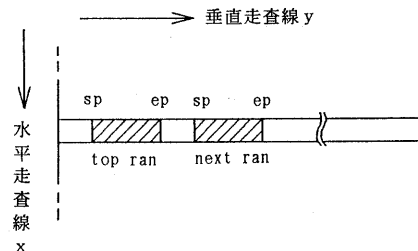
$$x_s \leq x \leq x_e$$

を満たす水平走査線 x に対して以下の処理を繰り返すことによって求められる。水平走査線 x 上の最初の横ランデータのポイント $h_{top}(x)$ により、水平走査線 x 上の最初の横ランデータを参照する。この横ランデータの始点座標 sp 、終点座標 ep が

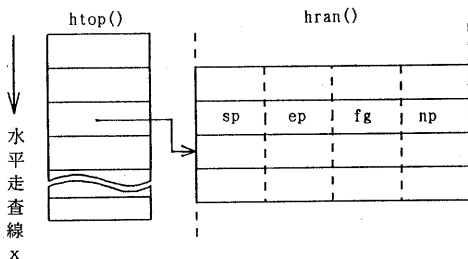
$$sp \leq y \leq ep$$

の条件を満たしているとき、この横ランが元の縦ランと交差するランである。そうでないときは、ポイント(np)により次のランを調べる。この操作を交差するランが求められるまで続ける。交差する横ランは1つの水平走査線に1つだけ存在するので、交差するランが見つかったら、次の水平走査線の処理に移る。

横ランと交差する縦ランも、同様にして効率よく求めることができる。



(a) 2値画像の水平走査線データ



(b) 横ランデータ構造

図5

4. 2 解析アルゴリズム

MCR表現を求める解析アルゴリズムのフローチャートを図6に示し、以下で説明する。

- ① 2値画像データを水平走査線ごとに読み込みながら縦ランと横ランのデータを登録し、同時に矩形解析も行う。この処理の詳しい説明は以下で示す。
- ② 未処理ランがなくなるまで、③と④の処理を繰り返す。
- ③ ②で検索された未処理ランと交差する未処理ランを求め、さらにそのランと交差する未処理ランを求めるといように、再帰的に交差する未処理ランを求めながら、隣接行列を作成する。
- ④ ③で作成した隣接行列に対して最大マッチング問題を解くことにより、最小被覆ランを求める。本研究ではHopcroft-Karpのアルゴリズムを用いて最大マッチングを解く(6)。

①のアルゴリズムについて詳しく説明する。ランデータの登録のために、現時点での水平走査線データと、1つ前の水平走査線データのためのラインバッファnewline()とoldline()を用意する。

横ランデータの登録は、ラインバッファnewline()の黒画素成分を調べることにより行う。

縦ランデータの登録は、2つのラインバッファを比較することにより行う。図7(a)に示すように、2つのラインバッファの比較により4つの状態変化

- I 白画素から白画素
・縦ランが存在していない
- II 白画素から黒画素
・縦ランが生成
- III 黒画素から黒画素
・縦ランが継続
- IV 黒画素から白画素
・縦ランが終了

が考えられる。状態変化IIでは、新しい縦ランを登録し始点座標を与える。状態変化IVでは、すでに登録してある縦ランの終点座標を与える。

図4(a)の形状の矩形解析について述べる。このような形状が存在するのは、長さmの縦ランがm個以上横方向に連続して終了した場合である。このことから状態変化IVの縦ラン終了時に、縦ランの長さと繰り返し回数を調べることにより、図4(a)の矩形領域を判断することができる。

次に図4(b)の矩形解析について述べる。この矩形領域が存在するのは、白画素から黒画素への境界と黒画素から白画素への境界のペア

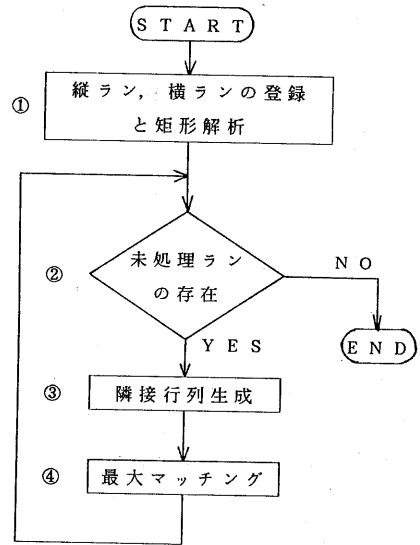
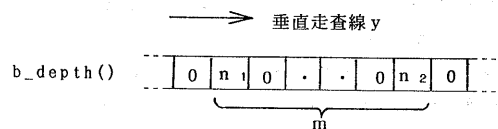
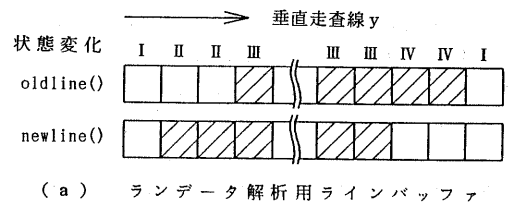


図6 MCR表現の解析アルゴリズム



(b) 図4(b)の矩形解析

図7 解析用データ

が幅 m をもち、 m 個以上縦方向へ続いた場合である。このため、境界の縦方向の連続数を示す配列 $b_depth()$ を用意する。 $b_depth()$ には白画素から黒画素への境界の縦方向の連続数、あるいは黒画素から白画素への境界の縦方向の連続数を設定し、境界でないときには 0 にする。 $b_depth()$ は水平走査線を読み込むごとに、常に更新する。図 7 (b) のように、 $b_depth()$ の境界の対 n_1 と n_2 、その間の幅 m に対して、

$$m \leq \min \{ n_1, n_2 \}$$

が成り立つとき、この領域が図 4 (b) の矩形領域にあたる。

図 4 (a)、(b) の矩形領域が検出された場合、該当するランを被覆ランか、非被覆ランに確定する。

5 実験結果

CCITT 標準原稿の MCR 表現を求め、その実験結果を示す。

CCITT 標準原稿は、大きさが縦 2376、横 1680 の 2 値画像データである。その例を図 8 (a) に示す。この例について MCR 表現を求め、その縦ラン被覆部、横ラン被覆部をそれぞれ図 8 (b)、(c) に示す。この実験結果により文書画像の表を構成する縦ストローク、横ストロークが、それぞれ縦ラン被覆部、横ラン被覆部のセグメントとして抽出されていることが分かる。文字に対してもストロークのはっきりしている部分は、ある程度ストロークと縦ラン被覆部、横ラン被覆部のセグメントとの対応がとれていることが分かる。

次に MCR 表現の応用例として、文書画像から表を抽出する問題を考えてみる。表を構成している各ストロークの長さは、文字の大きさに対して十分に長い。このことから適当な閾値を定めて、縦ラン被覆部、横ラン被覆部の各セグメントの形状を分析することにより表部分の抽出を行うことができる。図 8 (a) の文書画像から表を抽出した結果を図 8 (d) に示す。

6 まとめ

本研究では、2 値画像の新しい表現法として MCR 表現を提案した。この表現法は、2 値画像を最小個数の縦ランと横ランで被覆することにより画像を表現する方法である。2 値画像か

ら MCR 表現を求める問題は、縦ランと横ランを要素とする 2 部グラフの最大マッチング問題に帰着することを示した。さらに MCR 表現を求めるための、効率のよい解析アルゴリズムを明らかにした。また、文書画像に対する MCR 表現の実験結果を示し、文書画像解析への可能性を示した。

MCR 表現はできるだけ少ないランで 2 値画像を表現するという意味で情報圧縮を行っており、効率のよい文書画像の符号化の可能性をもっている。しかし、そのためには本手法に適合した符号化法を考える必要があり、今後の課題である。また、具体的な文書画像の解析問題に本表現法を適用して、その有効性を確かめることも課題として残されている。

謝辞 有益な御助言をいただいた東京工業大学精密工学研究所河原田弘教授に深謝いたします。

文献

- (1) "Open Document Architecture (ODA) and Interchange Format", CCITT, Rec. T.410 シリーズ (1988)
- (2) "Office Document Architecture (ODA) and Interchange Format", ISO 8613 (1988)
- (3) 山田, 藤長, 遠藤, 蓮池: "ミクスモード文書作成システム", 画電学誌, 15, 4 pp. 274-282 (1986)
- (4) 山田, 宮里, 蓮池: "マルチメディア文書構造処理システム", 画電学誌, 19, 5 pp. 286-295 (1990)
- (5) König, D.: "Graphen und Matrizen". Mat. Fiz. Lapok, 38, pp. 116-119. (1931)
- (6) Hopcroft, J.E. and Karp, R.M.: "An $n^{5/2}$ algorithm for maximum matchings in bipartite graphs", SIAM, J. Computer, 2, 4, pp. 225-231 (1973)
- (7) 嶋, 柏岡, 東野: "ランに対する座標演算に基づく 2 値画像の高速化回転のための一手法", 信学論 (D), J71-D, 7, pp. 1296-1305 (1988)

ÉTABLISSEMENT AJOURNÉ DÉPARTEMENT DU CAPITAL DE BIEN-ÊTRE 10, rue de la République, 7ème arrondissement Paris 7ème Téléphone: 270.11.00 Telex: 270.11.00 Fax: 270.11.00	Net Directeur CLASSEMENT FACTURE ÉMISSION 27011000 27011000 27011000 27011000 27011000 27011000	COMPTES 13 27011000 27011000 27011000 27011000 27011000 27011000
---	--	--

ÉTABLISSEMENT AJOURNÉ DÉPARTEMENT DU CAPITAL DE BIEN-ÊTRE 10, rue de la République, 7ème arrondissement Paris 7ème Téléphone: 270.11.00 Telex: 270.11.00 Fax: 270.11.00	Net Directeur CLASSEMENT FACTURE ÉMISSION 27011000 27011000 27011000 27011000 27011000 27011000	COMPTES 13 27011000 27011000 27011000 27011000 27011000 27011000
---	--	--

DATE SAISON	CODE BUDGET	COMPTES CLASSE	TYPE D'ORDRE	DATE DE LIQUIDATION	DATE DE DÉPART
01/01/77	12	ABCD	SP	15	99000

QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION
2	AP-809	Circuit intégré	2	104,35	F	208,66	F				
10	80-74	Composant	10	83,10	F	831,00	F				
25	8107	Composant indifférent	25	15,00	F	375,00	F				

DATE SAISON	CODE BUDGET	COMPTES CLASSE	TYPE D'ORDRE	DATE DE LIQUIDATION	DATE DE DÉPART
01/01/77	12	ABCD	SP	15	99000

QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION
2	AP-809	Circuit intégré	2	104,35	F	208,66	F				
10	80-74	Composant	10	83,10	F	831,00	F				
25	8107	Composant indifférent	25	15,00	F	375,00	F				

(a) 原画像 (C C I T T 標準原稿 NO 3)

(b) 縦ラン被覆部

DATE SAISON	CODE BUDGET	COMPTES CLASSE	TYPE D'ORDRE	DATE DE LIQUIDATION	DATE DE DÉPART
01/01/77	12	ABCD	SP	15	99000

QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION
2	AP-809	Circuit intégré	2	104,35	F	208,66	F				
10	80-74	Composant	10	83,10	F	831,00	F				
25	8107	Composant indifférent	25	15,00	F	375,00	F				

DATE SAISON	CODE BUDGET	COMPTES CLASSE	TYPE D'ORDRE	DATE DE LIQUIDATION	DATE DE DÉPART
01/01/77	12	ABCD	SP	15	99000

QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION	QUANTITÉ	UNITÉ	DESCRIPTION
2	AP-809	Circuit intégré	2	104,35	F	208,66	F				
10	80-74	Composant	10	83,10	F	831,00	F				
25	8107	Composant indifférent	25	15,00	F	375,00	F				

(c) 横ラン被覆部

(d) 表の部分抽出結果

図 8 文書画像の MCR 表現