# An Accuracy Barrier for Stable Three-Time-Level Difference Schemes for Hyperbolic Equations

R. Jeltsch, R. A. Renaut[1] and J. H. Smit [2]

[1]Department of Mathematics, Arizona State University, Tempe, AZ 85287-1804, USA
[2]Department of Mathematics, University of Stellenbosch, Stellenbosch 7600, South Africa

# An Accuracy Barrier for Stable Three-Time-Level Difference Schemes for Hyperbolic Equations

R. Jeltsch, R. A. Renaut[1] and J. H. Smit [2]

Seminar für Angewandte Mathematik
Eidgenössische Technische Hochschule
CH-8092 Zürich
Switzerland

## Abstract

We consider three-time-level difference schemes for the linear constant coefficient advection equation $u_t = cu_x$. In 1985 it was conjectured that the barrier to the local order $p$ of schemes which are stable is given by

$$p \leq 2\min\{R, S\}.$$

Here $R$ and $S$ denote the number of downwind and upwind points, respectively, in the difference stencil with respect to the characteristic of the differential equation through the update point. Here we prove the conjecture for a class of explicit and implicit schemes of maximal accuracy. In order to prove this result, the existing theory on order stars has to be generalized to the extent where it is applicable to an order star on the Riemann surface of the algebraic function associated with a difference scheme. Proof of the conjecture for all schemes relies on an additional conjecture about the geometry of the order star.

**Keywords:** scalar advection equation, difference scheme, accuracy, stability, order star, algebraic function, Riemann surface

**Subject Classification:** 65M10

---

[1]Department of Mathematics, Arizona State University, Tempe, AZ 85287-1804, USA

[2]Department of Mathematics, University of Stellenbosch, Stellenbosch 7600, South Africa

# 1  Introduction

Suppose we have a difference scheme for an initial-boundary value problem for a system of hyperbolic partial differential equations. A **global difference scheme** for the solution of this problem generally consists of an **interior scheme** and a **boundary scheme**. By the Lax-Richtmyer Equivalence Theorem these difference schemes result in solutions which are convergent to the exact solution only when they are consistent with the initial-boundary value problem and are stable. Consistency is the minimal requirement that the order of accuracy $p$ is one for interior and boundary schemes. Their stability in the global framework was investigated by Kreiss [18], [19] and in the influential paper by Gustafsson, Kreiss and Sundström [4]. In the latter paper the following necessary condition was given for the global scheme to be stable, namely that the corresponding interior scheme has to be stable in the Von Neumann sense when applied to the pure Cauchy problem for the scalar advection equation. Goldberg and Tadmor [3], [2] gave more practical sufficient conditions for stability of global schemes. These conditions entail, among others, that the boundary scheme also has to be stable in the same sense as mentioned above for the interior scheme.

From these results it can be concluded that accurate and stable difference schemes for the scalar advection equation are of fundamental importance in the construction of useful global schemes. For this reason we consider a Cauchy problem for the scalar advection equation

$$(1) \qquad\qquad \frac{\partial}{\partial t}\, u(t,x) \;=\; c\, \frac{\partial}{\partial x}\, u(t,x),\ x \in \mathbb{R},\ t \geq 0\ ,$$

$$u(0,x) \;=\; u_0(x)\ \text{given}\ ,$$

and a class of multistep $((k+1)$-time-level) difference schemes of the form

$$(2) \qquad\qquad \sum_{i=0}^{k} \sum_{j=-r_i}^{s_i} a_{ij}\, u_{n+i,m+j} = 0$$

which are used to determine an approximate solution of (1). The coefficients $a_{ij}$ are responsible for the two above-mentioned features, namely the **accuracy** and **stability** of the scheme. In general the requirement of stability imposes a bound on the order of a scheme. This paper focuses on this barrier to the order imposed by the requirement of stability for schemes of type (2).

One-step schemes $(k = 1)$ were extensively studied in [6, 7, 9, 10, 15, 22] and results for multistep schemes were given in [11, 13, 17, 23]. In [11, 15] it was conjectured that the order barrier for stable multi-time-level schemes (2) should be

$$(3) \qquad\qquad p \leq 2\,\min\{R, S\}\ .$$

Here, for counting purposes, we let the "zero line" be the characteristic through the point on the new time level for which one solves. Then $R$ denotes the number of downwind points

and $S$ the number of upwind points of a given scheme with respect to this zero line. This means that a stable scheme of order $p$ needs to have on each side of the characteristic at least $\lceil \frac{p}{2} \rceil$ points in the stencil. (Here $\lceil \alpha \rceil$ denotes the smallest integer which is not smaller than $\alpha$). If $p = 1$, this conjecture reduces to the Courant-Friedrichs-Lewy condition. Hence (3) has the quality of being an extension of the Courant-Friedrichs-Lewy condition, [1].

The bound (3) was proved in [15] for two-time-level schemes. In [17] it was partially proved for a small subclass of explicit three-time-level schemes. In [11], [13] many examples in support of (3) were given for multi-time-level schemes. In [12] the first lower bound in (3) for the $(k+1)$-level case was given by actually showing stability of schemes with long and slender stencils (only one step in space). Such schemes may be useful as high-order boundary schemes. In Sections 3-8 of this paper we generalize the results in [17] for convex maximal order explicit and implicit three-time level schemes (see Section 2). The results for all other schemes follow from a conjecture presented in Section 6.

The analysis in this paper is based on the order star technique, which was introduced in [24] and treated extensively in [5], [8]. These ideas have to be generalized for order stars on a Riemann surface defined by an algebraic function. This algebraic function is treated in Section 4. An additional complication is that our order stars are defined with respect to the comparison function $z^\mu$. This comparison function was first used in [22], [6]. The analysis in Sections 3-8 is a continuation of the work in [16], [17], viz. a study of the order stars on a two-sheeted Riemann surface. Since $z^\mu$ is multiple-valued with a logarithmic singularity at $z = 0$, extreme care has to be taken with the integration path used for the application of the argument principle. Notwithstanding these complications the order stars basically retain the elegant features which make them so useful in the sense that they allow a simple geometrical interpretation of the relationship between accuracy and stability.

For explicit schemes the order of the logarithmic singularity at $z = 0$ determines the maximum multiplicity of components of the order star. The various possible geometric configurations and the corresponding multiplicities of these geometries are investigated in Section 6. Except for a small subclass of schemes the derived bounds on order of the schemes do not lead to a proof of (3). This leads to the introduction of a conjecture that certain geometric configurations are not possible. A proof of the conjecture is provided for a subset of schemes of maximal order, see Section 7.

For implicit schemes the poles of the algebraic function also play an important role. The geometry of components containing poles is investigated in Section 8. Section 9 combines the results of the previous sections to provide the proof of (3).

We believe that this paper indicates the direction which the generalization of (3) to the $(k+1)$-time-level case will take. Since we restrict ourselves to schemes which can be considered convex, i.e. with an increasing stencil, we only work with poles of the algebraic function. In order to allow convexity for negative time as well, i.e. convexity in the reverse time direction, equivalently, concave schemes, the zeros of the algebraic

function also have to be taken into account. Note that by excluding concave schemes we exclude the possibility of a branch point of the algebraic function added to the logarithmic singularity at $z = 0$.

In a parallel investigation in [14] concerning three-time-level schemes for the wave equation we build on work started in [20], [21]. In that case the symmetry-properties of the schemes lead to a considerable simplification of the order star theory as treated in Sections 5-9. By also taking into account the role of the zeros of all of the polynomials defining the algebraic function the class of schemes can be treated there without imposing a restriction such as convexity. Furthermore, because of symmetry there is no possibility of a branch point at $z = 0$.

## 2  Order, stability and normalization of schemes

We consider three-time-level difference schemes of the form

$$(4) \qquad \sum_{j=-r_2}^{s_2} a_{2j}\, u_{n+2,m+j} + \sum_{j=-r_1}^{s_1} a_{1j}\, u_{n+1,m+j} + \sum_{j=-r_0}^{s_0} a_{0j}\, u_{n,m+j} = 0$$

$$n = 0, 1, 2, ..., \quad m = 0, \pm 1, \pm 2, ...$$

The step sizes in the time and space variables are denoted by $\Delta t$ and $\Delta x$, resp., while $\mu = \frac{c\Delta t}{\Delta x}$ denotes the Courant number which is assumed to be fixed. The coefficients $a_{ij}$ are real and depend in general on $\mu$, $a_{ij} = a_{ij}(\mu)$. Further $r_i$, $s_i \in \mathbb{Z}$ with $r_2 \geq 0$, $s_2 \geq 0$ and $-r_i \leq s_i$, $i = 0, 1$ and $a_{i,-r_i} \neq 0$, $a_{i,s_i} \neq 0$ for $i = 0, 1, 2$. The value $u_{nm}$ approximates $u(n\Delta t, m\Delta x)$. If $r_2 = s_2 = 0$ a scheme (4) is said to be **explicit**. Otherwise it is called **implicit**.

A scheme with a stencil satisfying

$$(5) \qquad \begin{cases} 0 & \leq & r_0 - r_1 \leq r_1 - r_2 \\ 0 & \leq & s_0 - s_1 \leq s_1 - s_2 \ . \end{cases}$$

is called a **convex scheme** with an **increasing stencil**. From a computational point of view these schemes seem to yield the most interesting stencils.

A Fourier Transform enables us to associate with (4) on time level $n + i$ a function

$$(6) \qquad a_i(z) = \sum_{j=-r_i}^{s_i} a_{ij}\, z^j, \qquad i = 0, 1, 2$$

and to introduce the **characteristic function**

$$\Phi(z, w) = a_2(z)\, w^2 + a_1(z)w + a_0(z) \ ,$$

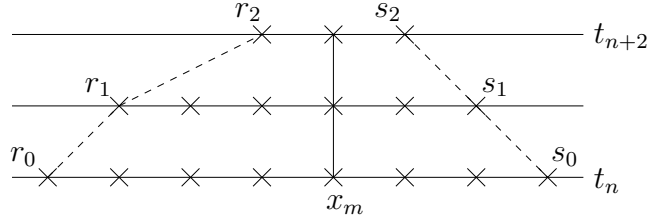which is assumed to be irreducible (see [13]).

3

Figure 1: Convex stencil

In order to be able to solve (4) for the values on the new time level in the implicit case, we impose the necessary and sufficient condition

$$a_2(z) \neq 0 \quad \text{for } |z| = 1 .$$

We also require our schemes to satisfy the following **normalization condition** (see [9, 13]):

(7) $$\begin{cases} r_2 &= \text{ number of zeros of } a_2(z) \text{ with } |z| < 1 \\ s_2 &= \text{ number of zeros of } a_2(z) \text{ with } |z| > 1 . \end{cases}$$

A scheme (4) is said to be **stable** if

(8) $$\left. \begin{array}{rcl} \Phi(z,w) &=& 0 \\ |z| &=& 1 \end{array} \right\} \implies \begin{cases} |w| \leq 1 \text{ and if } |w| = 1, \\ \text{then } w \text{ is a simple root} . \end{cases}$$

A scheme (4) has **error order** $p$ if for any smooth solution $u(t,x)$ of (1) we have

$$\sum_{i=0}^{2} \sum_{j=-r_i}^{s_i} a_{ij}\, u(t + i\Delta t,\; x + j\Delta x) \;\;=\;\; C\, \frac{\partial^{p+1}}{\partial x^{p+1}}\, u(t,x)(\Delta x)^{p+1} + O((\Delta x)^{p+2})$$

$$\text{if } \Delta x \to 0 \text{ and } \mu = \text{constant} .$$

Since we are interested only in schemes with positive order, we assume that

$$\Phi(1,1) = \sum_{i=0}^{2} \sum_{j=-r_i}^{s_i} a_{ij} = 0 .$$

The next result expresses the order of a scheme as a property of the solution $w$ of $\Phi(z,w) = 0$.

**Proposition 2.1** (Equivalent Order Conditions, [13, 23])**.** *Let a scheme (4) with characteristic function $\Phi(z,w)$ and Courant number $\mu$ be stable and satisfy $\Phi(1,1) = 0$. Then the following three conditions are equivalent.*

*a) The scheme has order $p$.*

4

*b)* $\Phi(z, z^\mu) = O((z-1)^{p+1})$ *as $z \to 1$.*

*c)* *The algebraic function $w$ given by $\Phi(z, w(z)) \equiv 0$ has exactly one branch $w_1$ which is analytic in a neighbourhood of $z = 1$ and satisfies*

$$z^\mu - w_1(z) = O((z-1)^{p+1}) \text{ as } z \to 1 .$$ $\qquad \square$

The next theorem gives the highest possible order that a scheme can have if stability is ignored. We introduce the index set of the difference stencil

$$I = \{(i,j) \in \mathbb{Z} \times \mathbb{Z} : 0 \le i \le 2, -r_i \le j \le s_i\} .$$

A scheme (4) is said to be **regular** if a characteristic line through any given stencil point does not pass through any other point of the difference stencil (see [13]).

**Proposition 2.2** (Regular Stencil, [13]). *Let a scheme (4) have a regular difference stencil with index set $I$. Then the highest possible order that the scheme can have is*

$$p = |I| - 2 ,$$

*where $|I|$ denotes the number of indices in $I$.* $\qquad \square$

# 3   Main result: bound on order of stable schemes

Suppose we have a convex scheme (4) which is also regular. The characteristic through the point $(t_{n+2}, x_m)$ (of the normalized scheme (4)) will be taken as the zero line. Then it is possible to interpret the order bound of stable schemes in a simple geometrical way such that for the highest order the number of points on each side of the zero line is balanced. If $R$ denotes the total number of downwind points and $S$ the total number of upwind points with respect to the zero line, then the order $p$ of a stable scheme satisfies

$$p \le 2 \min\{R, S\} .$$

This result can be related to the indices $r_i, s_i$ of (4) in the following way: Define

$$R_1 = \begin{cases} 0 & \text{if} \quad \mu < -r_1 \\ \lfloor r_1 + \mu \rfloor + 1 & \text{if} \quad -r_1 < \mu < s_1 , \\ r_1 + s_1 + 1 & \text{if} \quad \mu > s_1 \end{cases} \quad S_1 = r_1 + s_1 + 1 - R_1$$

$$R_0 = \begin{cases} 0 & \text{if} \quad 2\mu < -r_0 \\ \lfloor r_0 + 2\mu \rfloor + 1 & \text{if} \quad -r_0 < 2\mu < s_0 , \\ r_0 + s_0 + 1 & \text{if} \quad 2\mu > s_0 \end{cases} \quad S_0 = r_0 + s_0 + 1 - R_0 ,$$

where $\lfloor \alpha \rfloor$ denotes the largest integer not exceeding $\alpha$, and

(9) $$R = R_0 + R_1 + r_2, \quad S = S_0 + S_1 + s_2 .$$

Then the main result is as follows.

**Theorem 3.1** (Maximal Order of Stable Convex Schemes). *Let a convex scheme (4) with an increasing stencil be normalized and have a fixed Courant number $\mu$ satisfying $0 < |\mu| < \frac{1}{2}$. If the scheme is stable, then the order $p$ of the scheme is bounded by*

$$p \leq 2 \min\{R, S\} \ .$$

$\square$

**Remark 3.2.** *a) In [17] the bound (3) was proved for a small subclass of explicit schemes of type (4). In this paper we generalize it, making use of a conjecture introduced in Section 6, for the class of (explicit and implicit) schemes of type (4) which are convex and have an increasing stencil.*

*b) In Section 7 we provide a partial proof of (3). In particular, for the maximal order schemes, $p = |I| - 2$,*

$$p \leq \begin{cases} 2R & -\frac{1}{2} < \mu < 0 \\ 2S & 0 < \mu < \frac{1}{2}. \end{cases}$$

*c) The result (3) can be extended to $|\mu| > \frac{1}{2}$ by making use of the following transformation. Assume that a stable scheme is represented by $\Phi(z, w)$, where $w(z)$ approximates $z^\mu$ in a neighbourhood of the point $z = 1, w = 1$. Then we consider the scheme represented by the characteristic function*

$$\widetilde{\Phi}(z, u) = z^2 \, \Phi(z, \tfrac{u}{z}) \ .$$

*Since $u = zw$, the new scheme is stable and approximates $z^{\widetilde{\mu}} = z^{\mu+1}$ with the same order as the original scheme. The stencil undergoes the following transformations:*

$$\widetilde{r}_1 = r_1 - 1, \quad \widetilde{s}_1 = s_1 + 1, \quad \widetilde{r}_0 = r_0 - 2, \quad \widetilde{s}_0 = s_0 + 2 \ .$$

# 4 Properties of the algebraic function $w$

The algebraic function $w$, satisfying $\Phi(z, w(z)) \equiv 0$, is multiple-valued, consisting in general for a given $z$ of two values $w_1(z)$ and $w_2(z)$. Associated with this algebraic function is the Riemann surface $M$,

$$M = \{(z, w) \in \overline{\mathbb{C}} \times \overline{\mathbb{C}} : \ \Phi(z, w) = 0\} \ ,$$

consisting of two sheets, one above the other, interacting at a finite number of branch points $z_i$ (where $w_1(z_i) = w_2(z_i)$). The surface $M$ is a closed connected set on which $w$ is single-valued and, except for a finite number of singular points, also analytic.

**Remark 4.1** (Branch points of $w$). *a) The branch points of $w$ (except those that can occur at 0 and $\infty$) occur at points $z_i$ where $a_1^2(z_i) - 4a_2(z_i) \, a_0(z_i) = 0$. Since the coefficients of this polynomial equation are real, the branch points are either real or they occur in complex conjugate pairs. Branch cuts along which the two sheets of $M$ are connected can therefore always be taken to be straight lines which either fall on the real axis, or are orthogonal and symmetric to the real axes or occur in conjugate pairs.*

6

b) *If a scheme is stable, the corresponding algebraic function cannot have a branch point at $z = 1$ (see (8)). The sheet of $M$ on which the point $z = 1$, $w = 1$ occurs, is called the* **principal sheet**. *Since the Riemann surface is connected, this notion is basically a local property in a neighbourhood of $z = 1$. We make the convention that the principal sheet refers to that part of $M$ which can be connected to $z = 1$, $w = 1$ without crossing a branch cut. The remaining part of $M$ will be called the* **secondary sheet**.

**Remark 4.2** (Poles of $w$). *The function $w$ has a pole at every point where $a_2(z) = 0$. By the normalization condition (7) there are in total $r_2$ poles of $w$ with $|z| < 1$ away from $z = 0$ on the two sheets of $M$. Moreover, if $\max\{r_0, r_1, r_2\} > 0$, then $w$ can have a pole at $z = 0$ on one or both sheets of $M$ (see Proposition 4.4).*

**Remark 4.3** (Zeros of $w$). *The finite points where $w$ has zeros coincide with the points where the function $v(z) = 1/w(z)$ has poles. These points occur where $a_0(z)$ has zeros and occasionally also at $z = 0$.*

The expansion of $w(z)$ around $z = 0$, determined by use of Newton's polygons, is important in the subsequent discussion.

**Proposition 4.4** (Expansion at $z = 0$). *Let*

$$\Phi(z, w) = (a_{2, -r_2} z^{-r_2} + \ ... \ + a_{2, s_2} z^{s_2}) w^2 + (a_{1, -r_1} z^{-r_1} + \ ... \ + a_{1, s_1} z^{s_1}) w$$
$$+ (a_{0, -r_0} z^{-r_0} + \ ... \ + a_{0, s_0} z^{s_0})$$

*be the characteristic function of a convex scheme (4) with an increasing stencil. Then the algebraic function $w$ satisfying $\Phi(z, w) = 0$ does not have a branch point at $z = 0$ and has the following expansions at $z = 0$:*

a) *if $r_1 - r_2 > r_0 - r_1$, then*

$$w_1(z) \ = \ z^{-(r_1 - r_2)}(c_0 + c_1 z + c_2 z^2 + \ ...),$$
$$w_2(z) \ = \ z^{-(r_0 - r_1)}(d_0 + d_1 z + d_2 z^2 + \ ...).$$

*where $c_0 = -\dfrac{a_{1, -r_1}}{a_{2, -r_2}}$ and $d_0 = -\dfrac{a_{0, -r_0}}{a_{1, -r_1}}$.*

b) *if $r_1 - r_2 = r_0 - r_1$, then*

$$w_{1,2}(z) \ = \ z^{-(r_1 - r_2)}(-a_1(z)z^{r_1} \pm \sqrt{D})/(2a_2(z)z^{r_2})$$
$$D(z) \ = \ (a_1^2(z) - 4a_2(z)\, a_0(z))z^{2r_1}$$
$$= \ d + O(z)$$

*where $d = a_{1, -r_1}^2 - 4a_{2, -r_2}\, a_{0, -r_0}$.*

$\square$

**Remark 4.5.** *From Proposition 4.4 we observe that $z = 0$ is not a branch point of $w$ if $r_1 - r_2 > r_0 - r_1$. If $2r_1 = r_0 + r_2$ and $d \neq 0$ then again $z = 0$ is not a branch point. If $d = 0$ then $z = 0$ can be a branch point. In the following we shall restrict ourselves to schemes where $z = 0$ is not a branch point in which case one has the two expansions*

$$
\begin{aligned}
w_1(z) &= z^{-(r_1 - r_2)}(c_0 + c_1 z + c_2 z^2 + \dots) \\
w_2(z) &= z^{-(r_0 - r_1)}(d_0 + d_1 z + d_2 z^2 + \dots) .
\end{aligned}
$$

# 5   Order Stars

An order star is defined on the Riemann surface $M$ of the algebraic function $w$ in the following way. Define the function $\varphi$ by

$$
\varphi(z, w) = z^{-\mu} w, \quad (z, w) \in M
$$

and the **order star** $\Omega$ by

$$
\Omega = \{(z, w) \in M : \ |\varphi(z, w)| > 1\} .
$$

Because of the factor $z^{-\mu}$ the function $\varphi$ is multiple-valued on $M$. However, the order star $\Omega$, being defined by means of the modulus of $\varphi$, is again well defined on $M$. $\Omega^c$ denotes the complement of $\Omega$, i.e. $\Omega^c = M \backslash \Omega$. Because the coefficients $a_{ij}$ are real, $\Omega$ is symmetric with respect to the real axis.

The order and stability of a scheme, which were interpreted in Section 2 as properties of the function $w$, can be reinterpreted as properties of the order star. We give without proof those properties which are standard results in investigations involving order stars (e.g. in [24], [8]).

**Lemma 5.1** (Stability)**.** *If a scheme is stable, then*

$$
\Omega \cap \{(z, w) \in M : \ |z| = 1\} = \emptyset .
$$

$\square$

**Lemma 5.2** (Order)**.** *A scheme (4) has order $p$ if and only if at the point $z = 1$ on the principal sheet of $M$ the order star consists of $p + 1$ sectors of angle $\frac{\pi}{p+1}$, separated by $p + 1$ sectors of $\Omega^c$, each with the same angle.* $\square$

A subset $A$ (with boundary $\partial A$) of $\Omega$ is said to be an $\Omega$-**component** if $\partial A \subset \partial \Omega$ and $A$ is connected. $\Omega^c$-components are defined similarly. An $\Omega$-component is said to be of **multiplicity** $m$ if it contains $m$ $\Omega$-sectors at $z = 1$ on the principal sheet. Similarly for $\Omega^c$-components.

Note that the curve on $M$ which has the projection $|z| = 1$ in the $z$-plane separates $M$ into two well defined subsets. The set in $M$ with $|z| < 1$ is called the unit disk $\Delta$

and the set with $|z| > 1$ is called the outside of the unit disk. By Lemma 5.1 there is a clear distinction between the portion of the order star inside and the portion outside the unit disk. The components inside $\Delta$ are bounded, where a component $\Omega_1$ is said to be **bounded** if $\sup\limits_{(z,w)\in\Omega_1} |z| < \infty$.

In order to emphasize important features, our pictures of $\Omega$-components will not always be the exact geometrical embeddings of $M$ into $\mathbb{R}^3$. They will, however, display the basic connectivity relations and cuts, and elucidate the important properties of both macroscopic and microscopic scale.

According to Remark 4.5 we restrict ourselves to schemes where there is no branch point of $w$ at $z = 0$. Thus there are two values $w_1^0$ and $w_2^0$ with $\Phi(0, w_i^0) = 0$, i.e. there are two zero points $(0, w_1^0)$ and $(0, w_2^0)$ on $M$. Depending on the values of the indices $r_i$ we know from Proposition 4.4 that there can be a pole of $w$ at one or both of the zero points. Further, because of the factor $z^{-\mu}$ occurring in $\varphi$, this function in general no longer has an integer-valued leading exponent of $z$ at the zero points.

We know from Remark 4.2 that there are $r_2$ poles of $w$ away from $z = 0$ inside $\Delta$ (if our scheme is normalized). Since the expansion of $z^{-\mu}$ from any point $z_0 \neq 0$ has the form

$$z^{-\mu} = c_0 + c_1(z - z_0) + c_2(z - z_0)^2 + ..., \quad |z - z_0| < |z_0| \,,$$

there will be poles of $\varphi$ of exactly the same orders at the points with these $z$-values on one of the sheets of $M$.

The influence of poles and the behavior of $\varphi$ at $z = 0$ on the multiplicity of the components in which they occur, is studied by means of the argument principle with respect to the function $\varphi$ (see [24]). In this regard the factor $z^{-\mu}$ of $\varphi$ introduces onto $M$ a new structure in the sense that it defines on $M$ another Riemann surface which in general has infinitely many sheets. To define $z^{-\mu}$ uniquely on $M$, branch cuts, $L_i$, from $(0, w_i^0)$ to $(\infty, w_i^\infty)$, $i = 1, 2$ are made. These cuts are made according to the following rules.

**Rule 1 for cuts $L_i$:** The branch cuts $L_i$ have to be such that their projections onto the $z$-plane are either identical, or "enclose" a "sector" of $\mathbb{C}$ which does not contain a branch cut of $M$.

If we adhere to Rule 1, then $z^{-\mu}$ is defined uniquely on $M$ (see [17], Section 5.1), even if the cuts $L_i$ are allowed to cross a branch cut of $M$. In the present context, however, we can always avoid this. To this extent we introduce the following rule.

**Rule 2 for cuts $L_i$.** The cuts $L_i$ have to be such that each cut occurs only on one sheet of $M$, i.e. $L_1$ between $(0, w_1^0)$ and $(\infty, w_1^\infty)$ on the principal sheet and $L_2$ between $(0, w_2^0)$ and $(\infty, w_2^\infty)$ on the secondary sheet.

We can always adhere to Rules 1 and 2 by choosing the branch cuts $L_i$ to go along two
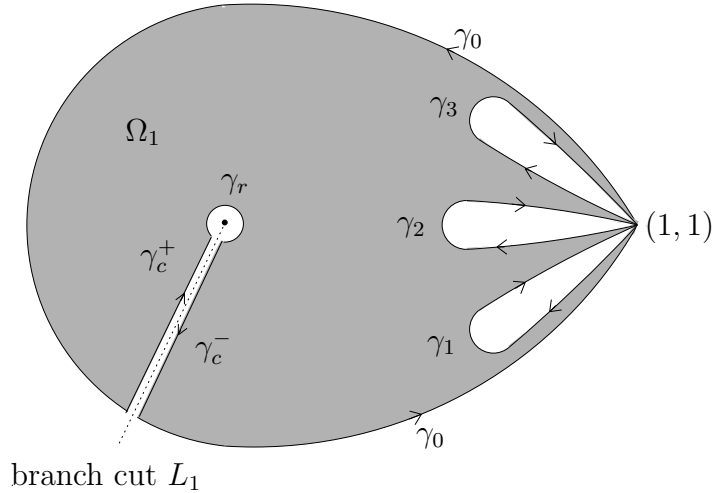
9

Figure 2: Component $\Omega_1$ illustrating the integration path

radial lines which have the same projection onto the $z$-plane and for which the projection onto the $z$-plane does not pass through the point $z = 1$. This convention for making cuts $L_i$ is adhered to unless explicitly indicated otherwise.

# 6    The Role of the Zero Points on Multiplicity

We start by restricting ourselves to the order stars of explicit schemes. Thus $r_2 = 0$ and there are no poles away from $z = 0$. Hence, inside $\Delta$, every bounded $\Omega$-component must contain at least one of the points $(0, \omega_i^0)$, $i = 1, 2$.

Our investigation of the relationship between the multiplicity of a component and the total order of poles/singularities of $\varphi$ that it contains begins with a very simple type of $\Omega$-component, $\Omega_1$ (say), which occurs only on the principal sheet of $M$ and which contains no branch points of $w$. This type of component was treated in [17], but the proof is repeated because it illustrates the appropriate application of the argument principle.

**Proposition 6.1** (Multiplicity). *Let $\Omega_1$ be such that the principal branch can be defined as a single-valued function on the projection of $\Omega_1$ onto the $z$-plane. Assume $\varphi$ has a leading exponent of $-\alpha$ at $z = 0$. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \leq \lfloor \alpha \rfloor + 1 \ .$$

**Proof:** If $\Omega_1$ is of multiplicity $m$, there are $m - 1$ $\Omega^c$-components emerging from (1,1) to the "inside" of $\Omega_1$. We evaluate

$$\frac{1}{2\pi i} \int_\gamma \frac{\varphi'(z, w)}{\varphi(z, w)} \, dz$$

10

where $\gamma$ is the closed curve which consists of the positively oriented boundary of $\Omega_1$ and a portion going around the zero point (see Fig. 2):

$\gamma_0$: positively oriented (w.r.t. $z = 0$) "outward" boundary of $\Omega_1$. According to [24] (proof of Proposition 4) the argument of $\varphi$ decreases along $\gamma_0$.

$\gamma_c^+, \gamma_c^-$: two sides of a Jordan curve which connects the "outward" boundary of $\Omega_1$ with a circle around $z = 0$. According to [17] (Lemma 4.4) the contributions to the integral along $\gamma_c^+$ and $\gamma_c^-$ cancel out.

$\gamma_r$: circular curve with small radius $r$, traversed clockwise. According to [17] (Lemma 4.3) the contribution of this curve to the integral is $\alpha$.

$\gamma_1, ..., \gamma_{m-1}$: boundary of $\Omega_1$ along $m - 1$ $\Omega^c$-components emerging from (1,1) to the "inside" of $\Omega_1$. Again the argument of $\varphi$ decreases along each path $\gamma_i$ and, because $\varphi$ is single valued in $\Omega_1 \backslash L_1$, every time the boundary $\gamma$ returns to (1,1) the argument has decreased by at least $2\pi$.

Then
$$\gamma = \gamma_0 + \gamma_c^+ + \gamma_r + \gamma_c^- + (\gamma_1 + \gamma_2 + ... + \gamma_{m-1}) .$$

By application of the argument principle, and because there are no zeros or poles of $\varphi$ inside $\gamma$, we have

$$
\begin{aligned}
0 &= \frac{1}{2\pi i} \int_\gamma \frac{\varphi'(z, w)}{\varphi(z, w)} \, dz \\
&= \underbrace{\frac{1}{2\pi i} \int_{\gamma_0}}_{<0} + \underbrace{\frac{1}{2\pi i} \left( \int_{\gamma_c^+} + \int_{\gamma_c^-} \right)}_{=0} + \underbrace{\frac{1}{2\pi i} \int_{\gamma_r}}_{=\alpha} + \underbrace{\frac{1}{2\pi i} \int_{\gamma_1 + ... + \gamma_{m-1}}}_{\leq -(m-1)} .
\end{aligned}
$$

Combining the first three terms and introducing the notation $\gamma_0^E$ to indicate the positively oriented curve
$$\gamma_0^E = \gamma_0 + \gamma_c^+ + \gamma_r + \gamma_c^- ,$$

it follows that
$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'(z, w)}{\varphi(z, w)} \, dz \leq \lfloor \alpha \rfloor .$$

Hence
$$m \leq \lfloor \alpha \rfloor + 1 . \qquad \square$$

Clearly the component being treated in Proposition 6.1 involves only one sheet of $M$, although nothing prevents it in general from crossing a branch cut of $M$ from one sheet to the other. We shall refer to a component of this kind as a **non-binary component**. With two sheets of $M$ available there also exist components which involve both sheets

of $M$ in a very specific manner and which will be called binary. We shall make these statements more precise in the following definition.

**Definition 6.2** (Binary/non-binary Components). *Let $\Omega_1$ be an $\Omega$-component containing exactly one zero point. Assume the branch cuts $L_i$ are radial lines with the same projection $L$ onto the z-plane and this projection does not pass through the point $z = 1$. We modify $\Omega_1$ into $\widetilde{\Omega}_1$ by making cuts along $L_i$ and encircling the zero points with infinitesimal small circles, see e.g. Fig. 2 and Fig. 4, such that $\widetilde{\Omega}_1$ satisfies the following properties*

   *i) $\widetilde{\Omega}_1$ is connected.*

   *ii) No closed curve in $\widetilde{\Omega}_1$ whose interior is contained completely in $\widetilde{\Omega}_1$ contains the zero point.*

   *iii) No projection of $\partial\widetilde{\Omega}_1$ onto the z-plane intersects with $L$.*

*$\Omega_1$ is called **non-binary** if the zero point is encircled once by such an infinitesimal small circle $\partial\widetilde{\Omega}_1$. In all other cases the component is called **binary**.*

If a component $\Omega_1$ has multiplicity $m$ its boundary $\partial\Omega_1$ can be decomposed naturally into $m$ curves $\gamma_i$ which connect the point $(1, 1)$. $\partial\widetilde{\Omega}_1$ consists also of $m$ curves $\widetilde{\gamma}_i$ which connect $(1, 1)$. These $\widetilde{\gamma}_i$ are either identical to $\gamma_i$ or are extended, $\gamma_i^E$, by a cut along $L_i$ and a circle around a zeropoint as was done with $\gamma_0$ in the previous proof.

**Lemma 6.3** (Symmetric Binary). *Let $\Omega_1$ be a symmetric binary component containing one zero point $(0, w_2^0)$ (say), with leading exponent $-\alpha_2$ of $\varphi$ at $(0, w_2^0)$, while the leading exponent of $\varphi$ at $(0, w_1^0) \notin \Omega_1$ is $-\alpha_1$. Let $\delta_i$ denote the non-integer part of $\alpha_i$, i.e.*

$$\delta_i = \alpha_i - \lfloor \alpha_i \rfloor, \;\; i = 1, 2 \; .$$

*Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$(10) \qquad\qquad m \leq 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor \; .$$

**Proof:** The proof is conducted in two different ways depending on the way in which the branch cuts $L_i$ are chosen. The first version highlights the binary character of the component, while the second leads to the bound (10).

**Version 1:** We first deviate from our convention of making the cuts $L_i$ by choosing them such that their projection onto the $z$-plane is the positive semi axis (see Fig. 3). With the integration along $\gamma_0^E$ we have the situation that, on the principal sheet, we have gone once around the zero point $(0, w_1^0)$ without crossing $L_1$. Hence we end up with $\gamma_0^E$ at a point where $z^{-\mu}$, and therefore also $\varphi$, has a value which differs from the value with which it started at the point $(1,1)$. We then have to return along $\gamma_1$ to our starting point at $(1,1)$. Only then has the complete boundary of a component with respect to the function
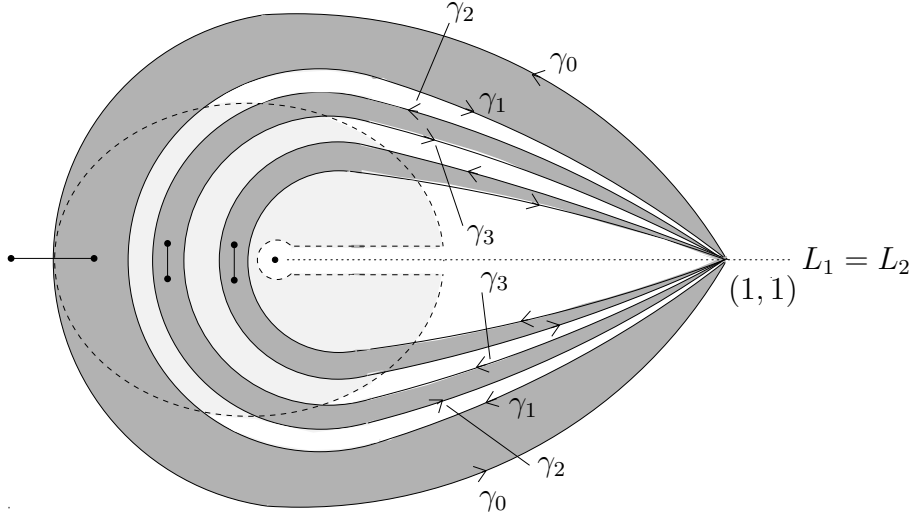
12

Figure 3: Symmetric binary component and cuts $L_i$ going through $z = 1$

$\varphi$ been traversed. Hence, in terms of counting the multiplicity of $\Omega_1$ at $(1,1)$, we can regard the point where $\gamma_0^E$ went over into $\gamma_1$ as a point away from $(1,1)$. When we apply the argument principle as we did in Proposition 6.1, this only accounts for the sectors of $\Omega_1$ "on one side" of the cut $L_1$. Disregarding the sectors in the lower halfplane $Im\ z < 0$ we have a component with $m/2$ sectors in the upper halfplane which is, with respect to these $m/2$ sectors, a non-binary component. Hence

$$\frac{1}{2\pi i} \int_{\gamma_0^E + \gamma_1} \frac{\varphi'(z, w)}{\varphi(z, w)}\, dz \le \lfloor \alpha_2 \rfloor$$

and as in the proof of Proposition 6.1,

(11) $$0 \le \lfloor \alpha_2 \rfloor - \left( \tfrac{m}{2} - 1 \right)$$

from which we obtain

$$m \le 2\lfloor \alpha_2 \rfloor + 2\ .$$

In this case the integration along $\gamma_2 + \gamma_3$, which has led to a drop in the argument by $2\pi$, has contributed two sectors as compared to just one for a non-binary component.

**Version 2:** In the version 1 proof the cuts $L_i$ were not made according to convention. By making the cuts according to convention (see Fig. 4) we find

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'(z, w)}{\varphi(z, w)}\, dz \le \lfloor \alpha_1 + \alpha_2 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_1^E} \frac{\varphi'(z, w)}{\varphi(z, w)}\, dz \le \lfloor -\alpha_1 \rfloor\ .$$

and

$$\frac{1}{2\pi i} \int_{\gamma_2^E} \frac{\varphi'(z, w)}{\varphi(z, w)}\, dz + \frac{1}{2\pi i} \int_{\gamma_3^E} \frac{\varphi'(z, w)}{\varphi(z, w)}\, dz \le \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor = -1\ .$$
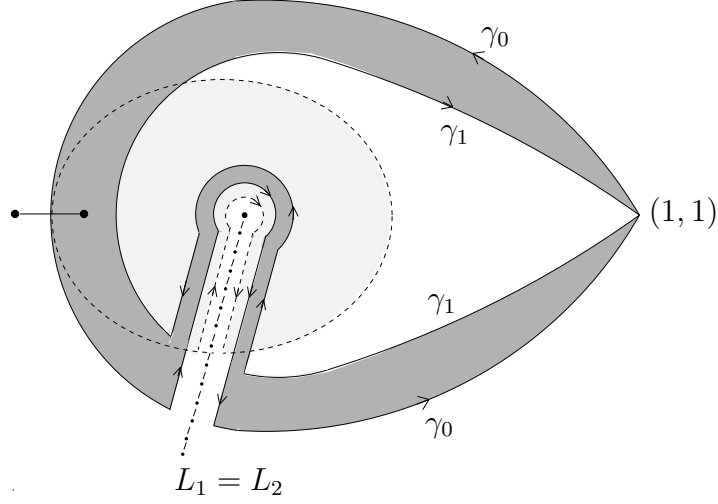
13

Figure 4: Symmetric binary component and cuts $L_i$ made according to convention

Assuming $m/2 - 1$ pairs of curve like $\gamma_2 + \gamma_3$ application of the argument principle yields

(12)
$$0 \leq \lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor - (\tfrac{m}{2} - 1) \,.$$

Hence

$$m \leq 2\{\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor) + \lfloor -\alpha_1 \rfloor\} + 2$$

and

$$m \leq 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor \,,$$

where we have made use of the fact that

$$\lfloor \alpha \rfloor + \lfloor -\alpha \rfloor = -1 \ \text{if} \ \alpha \notin \mathbb{Z} \,. \qquad \square$$

**Remark 6.4.** *a) The zero point not inside $\Omega_1$ can be excluded by traversal of more than one "inner" boundary curve which crosses the negative real axis: see Fig. 5.*

If the integration process is carried out in is the proof of Lemma 6.3, we obtain

$$\frac{1}{2\pi i} \ \int_{\gamma_0^E} \leq \lfloor \alpha_1 + \alpha_2 \rfloor, \ \ \frac{1}{2\pi i} \ \int_{\gamma_1^E} \leq \lfloor -\alpha_1 \rfloor, \ \ \frac{1}{2\pi i} \ \int_{\gamma_2^E} \leq \lfloor \alpha_1 \rfloor, \ \ \frac{1}{2\pi i} \ \int_{\gamma_3^E} \leq \lfloor -\alpha_1 \rfloor \,,$$

leading to the inequality

$$0 \leq \{\lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor + \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor\} - \frac{(m-4)}{2} \,,$$

with $m$ again bounded by

$$m \leq 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor \,.$$

*b) If $(0, w_1^0)$ is contained in $\Omega_1$ and $(0, w_2^0)$ is excluded the multiplicity of $\Omega_1$ is determined in the same way but with $\lfloor \alpha_2 \rfloor$ replaced by $\lfloor \alpha_1 \rfloor$, see Fig. 6.*
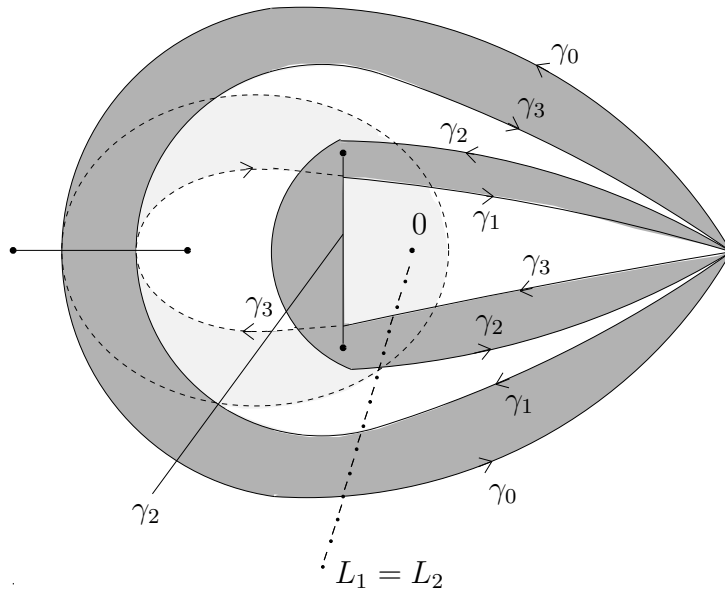
14

Figure 5: Symmetric binary component with three "inner" boundary curves
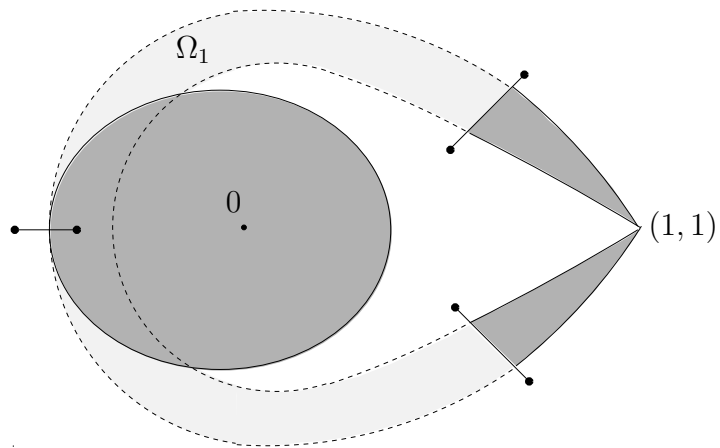


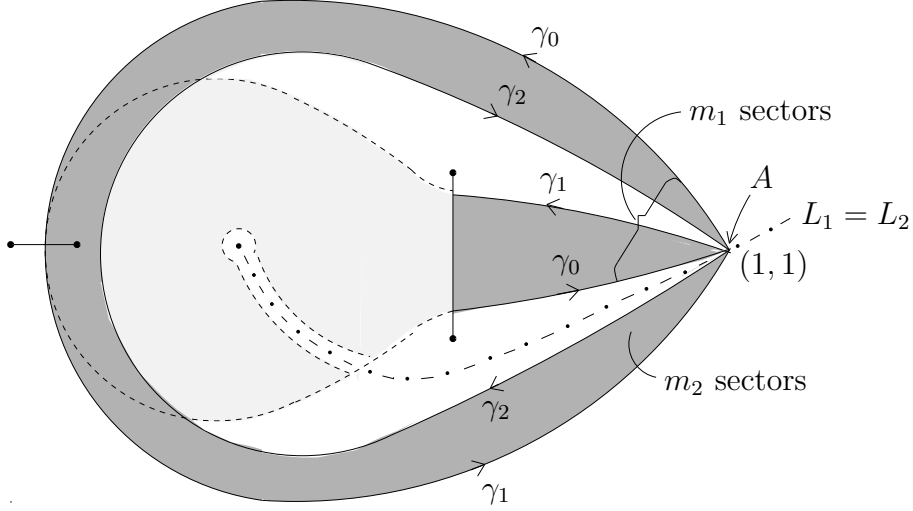Figure 6: Binary component containing the zero point on the principal sheet

Figure 7: Non-symmetric binary component and cuts $L_i$ going through $z = 1$

**Lemma 6.5** (Non-symmetric Binary). *Let $\Omega_1$ be a non-symmetric binary component containing one zero point $(0, w_2^0)$ (say), with leading exponent $-\alpha_2$ of $\varphi$ at $(0, w_2^0)$. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \leq 2\lfloor \alpha_2 \rfloor + 1 .$$

**Proof:** The proof is again conducted in two ways as a result of two different choices of the branch cuts $L_i$.

**Version 1:** If cuts $L_1$ and $L_2$ were chosen as in version 1 of the proof of Lemma 6.3 they would intersect a branch cut of $M$. Therefore the cuts $L_1$ and $L_2$ are made such that the projection winds from $(0, 0)$ to $(1, 1)$ without intersecting either the projection of $\partial\Omega$ onto $\mathbb{C}$ or the projection of any cut of $M$ onto $\mathbb{C}$, see Figure 7. Then $\Omega_1$ is non-symmetric with respect to $L_1$. The argument principle is applied along the positively oriented boundary $\gamma = \gamma_0^E + \gamma_1 + \gamma_2$ which starts out on one side of the cut $L_1$ and is assumed to consist of $m_1$ sectors at $(1, 1)$. The contributions of the integrals are

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'(z, w)}{\varphi(z, w)} \, dz \leq \lfloor \alpha_2 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_1 + \gamma_2} \frac{\varphi'(z, w)}{\varphi(z, w)} \, dz \leq -1 .$$

Application of the argument principle leads to

(13) $$0 \leq \lfloor \alpha_2 \rfloor - 1 - (m_1 - 2) .$$

Observing that in the optimal case $m_1 = (m - 1)/2 + 1$ we obtain the result.

The process is now repeated for the portion of $\Omega_1$, "on the other side" of $L_1$, where we assume a total of $m_2$ sectors. Integrating along $\gamma_2$ we end up at a point $A$ (say) where
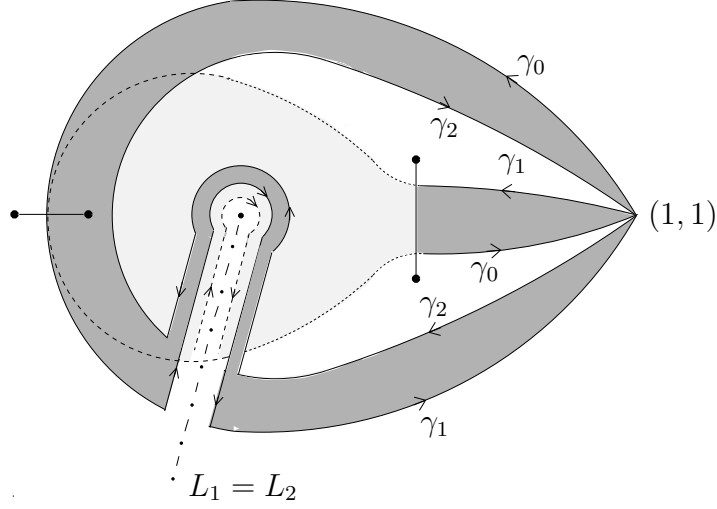
16

Figure 8: Non-symmetric binary component and cuts $L_i$ made according to convention

$\varphi$ is different to the initial value. Continuing along $\gamma_0^E$ to $A$ again, and then along $\gamma_1$ to return to our starting point we obtain

$$\frac{1}{2\pi i} \int_{\gamma_2 + \gamma_0^E + \gamma_1} \frac{\varphi'(z,w)}{\varphi(z,w)} \, dz \leq \lfloor \alpha_2 - 1 \rfloor \,,$$

where the $-1$ accounts for the fact that we returned to $A$, on which occasion the argument must have decreased by at least $2\pi$. Application of the argument principle leads to

(14) $$0 \leq \lfloor \alpha_2 \rfloor - 1 \, - \, (m_2 - 1) \,.$$

By combining (13 )and (14) we obtain the following bound on the total number $m = m_1 + m_2$ of sectors of $\Omega_1$:

$$m \leq 2\lfloor \alpha_2 \rfloor + 1 \,.$$

**Version 2:** By choosing the branch cuts $L_i$ according to convention (see Fig. 8), we obtain the following for the integrals:

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor \alpha_2 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_1^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor \alpha_1 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_2^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor -\alpha_1 \rfloor \,.$$

Application of the argument principle leads to

$$0 \leq \{\lfloor \alpha_2 \rfloor + \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor\} - \frac{(m-3)}{2} \,,$$

where the factor 2 accounts for the binary nature of the component and 3 is subtracted from $m$ because $\gamma_0, \gamma_1$ and $\gamma_2$ contribute 3 sectors. Hence, we again obtain
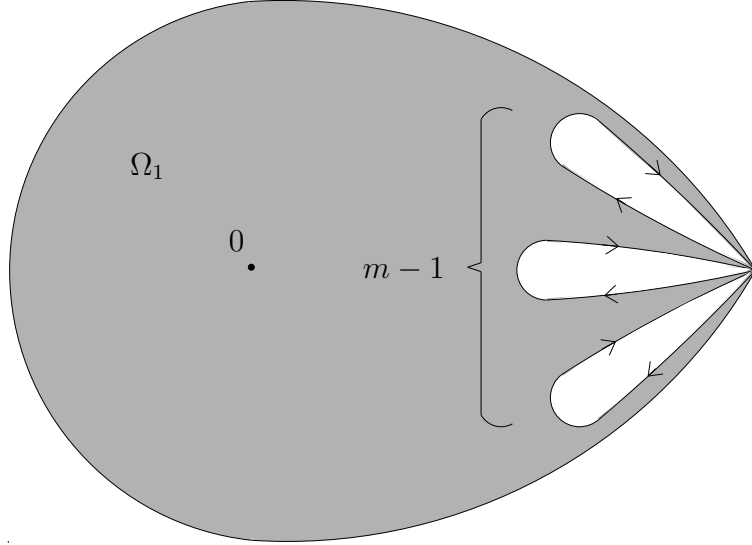
$$m \leq 2\lfloor \alpha_2 \rfloor + 1 \,.$$

$\square$

Figure 9: Non-binary component with $m-1$ $\Omega^c$-components

In view of Lemmas 6.3 and 6.5 we have to conclude that the bound (10) is in general too sharp if the non-integer parts of $\alpha_1$ and $\alpha_2$ satisfy $0 < \delta_1 + \delta_2 < 1$. A combination of Lemmas 6.3 and 6.5 leads to the following general result for binary components.

**Proposition 6.6** (Binary). *Let $\Omega_1$ be a binary component containing one zero point. Assume that $\varphi$ has a leading exponent of $-\alpha_1$ at the zero point inside $\Omega_1$ and $-\alpha_2$ at the other zero point. Let $\delta_1$ and $\delta_2$ be the non-integer parts of $\alpha_1$ and $\alpha_2$, respectively, i.e.*

$$\delta_1 = \alpha_1 - \lfloor \alpha_1 \rfloor, \quad \delta_2 = \alpha_2 - \lfloor \alpha_2 \rfloor .$$

*Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \le 2\lfloor \alpha_1 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\} .$$

$\square$

**Remark 6.7** (Efficiency). *a) In view of the factor $2$ accompanying $\lfloor \alpha \rfloor$ in the bound for binary components, we say that the zero point inside a binary component has a **higher efficiency** than the zero point inside a non-binary component. The zero point is in this case regarded as contributing twice to the multiplicity of the component.*

*b) The multiplicity of a non-binary component $\Omega_1$ is achieved by $m-1$ $\Omega^c$ components which are bounded by $\Omega_1$ and do not loop around either zero point, Figure 9. On the contrary, the multiplicity of a binary component is achieved because of $\Omega^c$ components which do loop around one of the zero points. In [17] these were referred to as binary loops, (Lemma 5.10). Moreover, the connectedness of $\Omega_1$ requires branch cuts.*
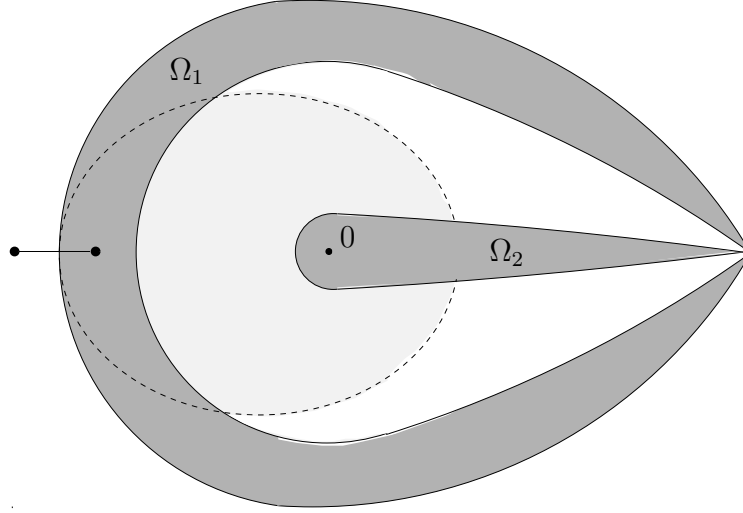
18

Figure 10: Binary component with one "binary loop"

**Proposition 6.8.** *There can be at most one binary component containing one zero point inside the unit disk $\Delta$.*

**Proof:** Let $\Omega_1$ be a binary component inside $\Delta$ and say $(0, w_2^0)$ is contained in $\Omega_1$, while $(0, w_1^0)$, which does not belong to $\Omega_1$, is enclosed by (an) "inner" boundary curve(s) of $\Omega_1$ (see Definition 6.2). Suppose $(0, w_1^0)$ belongs to a second $\Omega$-component $\Omega_2$ (say). For $\Omega_2$ to be binary, it has to have (an) "outward" boundary curve(s) going through a branch cut on the negative real axis and then enclosing $(0, w_1^0)$. However, this is impossible since $(0, w_1^0)$ is already enclosed by (an) "inner" boundary curve(s) of $\Omega_1$. Hence, $\Omega_2$ cannot be binary. □

A binary component $\Omega_1$ can be combined with a non-binary component $\Omega_2$ (say), as is illustrated in Fig. 11. For such a combination the following theorem follows as a consequence of Propositions 6.1 and 6.6.

**Theorem 6.9** (Binary plus non-binary). *Let $\varphi$ have leading exponents of $-\alpha_1$ and $-\alpha_2$ at the zero points $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$ and suppose $(0, w_1^0)$ belongs to a non-binary component and $(0, w_2^0)$ to a binary component. Then the highest total multiplicity $m$ that these two components can contribute at (1,1) is given by*

(15) $$m \leq (\lfloor \alpha_1 \rfloor + 1) + (2\lfloor \alpha_2 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\}),$$

*where $\delta_i = \alpha_i - \lfloor \alpha_i \rfloor$, $i = 1, 2$.* □

**Remark 6.10.** *The bound (15) is a sharper bound than (5.20) in [17], p. 29, Proposition 5.15, because the contribution due to $\lfloor \alpha_1 \rfloor$ is not doubled.*

19

Figure 11: Binary component $\Omega_1$ combined with non-binary component $\Omega_2$

## 6.1 Components containing two zero points

An obvious way of obtaining a component with two zero points is to connect two components with one zero point each by means of a branch cut of $M$. We shall show that our technique to prove bounds for the multiplicity in such a situation will give a bound which is larger than what one would obtain by ignoring the connecting cut and applying the results of the previous section to each component separately. Since we have not found any example which shows that this higher bound is sharp we conjecture that the smaller bound (15) is correct in all cases.

**Conjecture 6.11.** *Let $\varphi$ have leading exponents of $-\alpha_1$ and $-\alpha_2$ at the zero points $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$. Then the multiplicity $m$ of the $\Omega-$component $\Omega_1$ containing both zero points satisfies*

$$ m \leq \lfloor \alpha_1 \rfloor + 1 + 2 \lfloor \alpha_2 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\} $$

Clearly, if $\Omega_1$ in the conjecture can be separated into two components then the conjecture is proved.

We can prove this conjecture for a class of schemes of maximal order, $p = |I| - 2$. The proof is obtained by contradiction. Hence we assume the converse and examine its implications.

**Definition 6.12** (Double Binary Component). *An $\Omega$ component $\Omega_1$, which contains both zero points is called double binary if the contribution to its multiplicity by both zero points occurs via a doubling of $\lfloor \alpha_1 \rfloor$ and $\lfloor \alpha_2 \rfloor$.*

Figure 12: Symmetric double-binary component

A double-binary component can be either symmetric or non-symmetric. We first consider Figure 12 which depicts a modification of the symmetric binary component illustrated in Figs. 3, 4, in the sense that the "inner" boundary curve is moved inward so as to include the second zero point, and hence give a component containing both zero points.

**Lemma 6.13** (Symmetric Double Binary Component). *Let $\Omega_1$ be a symmetric double-binary component containing both zero points, $(0, w_1^0)$ and $(0, w_2^0)$ with leading exponents $-\alpha_1$ and $-\alpha_2$ of $\varphi$ at $(0, w_1^0)$ and $(0, w_2^0)$, respectively. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \le 2[\alpha_1] + 2[\alpha_2] + 1 + 2[\delta_1 + \delta_2].$$

**Proof:** Clearly if cut $A$ would not be present one could apply Theorem 6.9 and we would have Conjecture 6.11. But, since the cut $A$ can't be removed we apply the argument principle to the whole component $\Omega_1$. Hence

$$
\begin{aligned}
0 &= \frac{1}{2\pi i} \int \frac{\gamma'}{\gamma} \, dz \\
&= \frac{1}{2\pi i} \int_{\gamma_0^E} + \frac{1}{2\pi i} \int_{\gamma_1^E} + \frac{1}{2\pi i} \sum_{j=1}^{\frac{m-3}{2}} \left( \int_{\gamma_{2j}^E} + \int_{\gamma_{2j+1}^E} \right) + \frac{1}{2\pi i} \int_{\gamma_{m-1}^E} \\
&\le \lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor + \tfrac{m-3}{2}(-1) + \lfloor \alpha_1 \rfloor
\end{aligned}
$$

and this gives the bound

$$(16) \qquad m \le 2\lfloor \alpha_1 \rfloor + 1 + 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor .$$

$\square$

21

Figure 13: Non-symmetric double binary

Clearly, if $\lfloor \alpha_1 \rfloor > 0$ then the bound (16) is not as sharp as the conjecture bound. While in this example the component could have been separated by removing cut $A$ such a separation is not as evident in the example depicted in Fig. 13. A comparison with Figs. 7 and 8 reveals that Fig. 13 is obtained from a non-symmetric binary component by replacing the cut locally orthogonal to the real axis by a cut, cut $A$, lying on the projection of the real axis onto $C$. Therefore a component of this kind will be called non-symmetric double-binary.

**Lemma 6.14** (Non-symmetric Double Binary)**.** *Let $\Omega_1$ be a non-symmetric double-binary component containing both zero points, $(0, w_1^0)$ and $(0, w_2^0)$, with leading exponents $-\alpha_1$ and $-\alpha_2$ of $\varphi$ at $(0, w_1^0)$ and $(0, w_2^0)$, respectively. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \leq 2(\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1).$$

**Proof:** As usual we apply the argument principle to the whole component $\Omega_1$. Hence

$$
\begin{aligned}
0 &= \frac{1}{2\pi i} \int \frac{\gamma'}{\gamma} \, dz \\
&= \frac{1}{2\pi i} \int_{\gamma_0^E} + \frac{1}{2\pi i} \int_{\gamma_1^E} + \frac{1}{2\pi i} \int_{\gamma_2^E} + \frac{1}{2\pi i} \sum_{j=1}^{\frac{m-4}{2}} \left( \int_{\gamma_{2j+1}^E} + \int_{\gamma_{2j+2}^E} \right) + \frac{1}{2\pi i} \int_{\gamma_{m-1}^E} \\
&\leq \lfloor \alpha_2 \rfloor + \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor + \tfrac{m-4}{2} \left( \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor \right) + \lfloor \alpha_1 \rfloor
\end{aligned}
$$

and this gives the bound

$$(17) \qquad\qquad m \leq 2\lfloor \alpha_1 \rfloor + 1 + 2\lfloor \alpha_2 \rfloor + 1 \,. \qquad\qquad \Box$$

22

Again, if $\lfloor \alpha_1 \rfloor > 0$ then the bound (17) is not as sharp as the conjectured bound. As in the previous example the bound is wrong by the factor 2 in the term containing $\lfloor \alpha_1 \rfloor$.

We observe that the results suggested by Lemmas 6.13 and 6.14 lead to the following general result for a double binary component.

**Theorem 6.15** (Double Binary). *Let $\varphi$ have leading exponents of $-\alpha_1$ and $-\alpha_2$ at $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$, and suppose that both $(0, w_1^0)$ and $(0, w_2^0)$ belong to one double-binary component, $\Omega_1$. Then the highest multiplicity, $M$, that this component can contribute at $(1, 1)$ is given by*

$$m \le 2(\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor) + 1 + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\}$$

$\square$

# 7  The role of the branch cuts on multiplicity

Close inspection of the results derived until now will reveal that these results actually rely on the occurrence of branch cuts to connect binary loops to the portion of the component containing one or both zero points. If there are insufficient branch cuts, or equivalently insufficient branch points, these multiplicities will not be achievable. It therefore becomes appropriate to reformulate the earlier results in terms of the minimum number of branch points used by a component in order to achieve a certain multiplicity. The implication is that we now consider components which we will call **suboptimal**; for example, for a binary component, $\Omega_1$, we allow for the possibility that there are insufficient branch points for the zero point to be completely binary, and hence that there are also sectors of $\Omega_1$ at $z = 1$ which contribute only a factor 1 rather than a factor 2 to the multiplicity.

To standardize our approach we will adopt the following notation:

$$
\begin{aligned}
K &:= \text{ number of branch points utilized by the component} \\
m_1 &:= \text{ number of sectors at } z = 1 \text{ due to binary loops} \\
m_2 &:= \text{ number of sectors at } z = 1 \text{ due to non-binary loops.}
\end{aligned}
$$

We also make the assumption, without loss of generality, $\lfloor \alpha_2 \rfloor \ge \lfloor \alpha_1 \rfloor$. Furthermore, to ease comparison with the previous results, we will refer to the components as classified as in the earlier sections. In each case we will employ the "Version 2" type proofs since these give the tighter bounds.

**Lemma 7.1** (Suboptimal Symmetric Binary Component, (cf. Lemma 6.3)). *Let $\Omega_1$ be a symmetric binary component containing one zero point $(0, w_2^0)$, (say), with leading exponent $-\alpha_2$ of $\varphi$ at $(0, w_2^0)$, while the leading exponent of $\varphi$ at $(0, w_1^0) \notin \Omega_1$ is $-\alpha_1$. Further,*
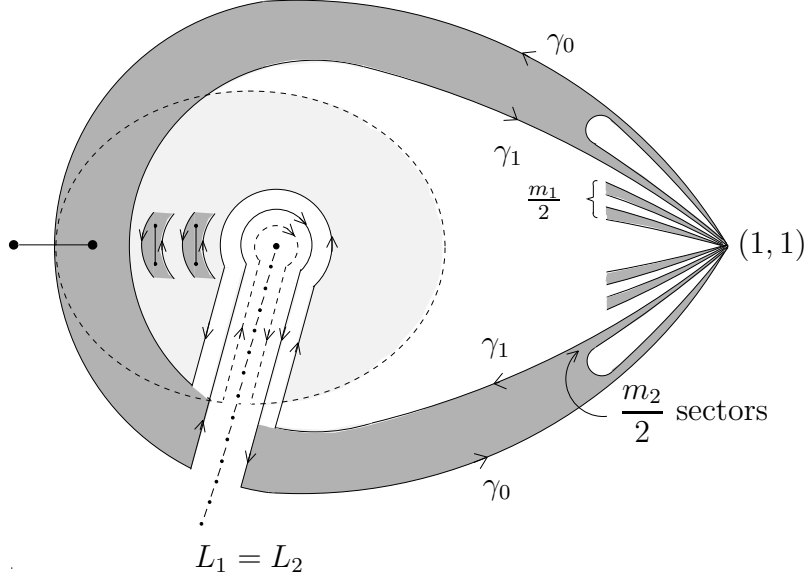
Figure 14: Suboptimal symmetric binary

*suppose that $\Omega_1$ contains at most $K$ branch points of $w(z, \mu)$. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \le \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor + \min\{\lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor, \lfloor \frac{K+1}{2} \rfloor\}.$$

**Proof:** The proof follows as for Lemma 6.3 but note now that the number of binary loops is limited by the number of branch points $K$. Each of the $\frac{m_1}{2}$ binary loops contributes $-1$ to the argument but two sectors at $z = 1$. The $m_2$ non-binary loops also contribute $-1$ to the argument decrease but one sector at $z = 1$. Hence, the total number of sectors at $z = 1$ is

$$m = m_1 + m_2 + 2.$$

Further, by the argument principle

$$0 \le \lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor - \frac{m_1}{2} - m_2.$$

Therefore

$$\begin{aligned}
m = m_1 + m_2 + 2 \quad &\le \quad \lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor + \frac{m_1}{2} + 2 \\
&= \quad \lfloor \alpha_2 \rfloor - 1 + \lfloor \delta_1 + \delta_2 \rfloor + 2 + \frac{m_1}{2}.
\end{aligned}$$

Now each binary loop uses at least one branch cut to connect that loop to $\Omega_1$. Also one branch point is required to make the component binary. Therefore

$$m_1 \le K - 1,$$

24

Figure 15: Suboptimal non-symmetric binary

and

$$\frac{m_1}{2} \leq \lfloor \frac{K-1}{2} \rfloor.$$

Thus

$$m \leq \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor + 1 + \lfloor \frac{K-1}{2} \rfloor$$

and by Lemma 6.3 the result follows. □

**Corollary 7.2.** *The minimum number of branch points, $K_{\min}$, contained in a component $\Omega_1$, described as in Lemma 7.1, for which the maximum multiplicity, as indicated by Lemma 6.3, is obtained, is given by*

$$\lfloor \frac{K_{\min}+1}{2} \rfloor = \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor.$$

**Proof:** This follows immediately by observing that when $m_2 = 0$, $K \geq m-1$, and hence $K_{\min} = m_{\max} - 1$, where $m_{\max} = 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor$. □

A non-symmetric binary component can also be suboptimal.

**Lemma 7.3** (Suboptimal Non-symmetric Binary Component, (cf. Lemma 6.5)). *Let $\Omega_1$ be a non-symmetric binary component containing one zero point $(0, w_2^0)$, (say), with leading exponent $-\alpha_2$ of $\varphi$ at $(0, w_2^0)$, while the leading exponent of $\varphi$ at $(0, w_1^0) \notin \Omega_1$ is*

25

$-\alpha_1$. *Further, suppose that* $\Omega_1$ *contains at most* $K$ *branch points of* $w(z, \mu)$. *Then the multiplicity* $m$ *of* $\Omega_1$ *satisfies*

$$m \leq \lfloor \alpha_2 \rfloor + \min\{\lfloor \alpha_2 \rfloor + 1, \lfloor \frac{K+1}{2} \rfloor\}.$$

$\square$

**Corollary 7.4.** *The minimum number of branch points,* $K_{\min}$, *contained in a component* $\Omega_1$, *described as in Lemma 7.3, for which the maximum multiplicity, as indicated by Lemma 6.5, is obtained, is given by*

$$\lfloor \frac{K_{\min}+1}{2} \rfloor = \lfloor \alpha_2 \rfloor + 1.$$

$\square$

These suboptimal binary components can be combined with a non-binary component in exactly the same way as a binary component is combined with a non-binary component in Theorem 6.9:

**Theorem 7.5** (Suboptimal Binary plus Non-binary). *Let* $\varphi$ *have leading exponents of* $-\alpha_1$ *and* $-\alpha_2$ *at the zero points* $(0, w_1^0)$ *and* $(0, w_2^0)$, *respectively, on the two sheets of* $M$, *and suppose* $(0, w_1^0)$ *belongs to a nonbinary component while* $(0, w_2^0)$ *belongs to a binary component containing* $K$ *branch points. Then the highest total multiplicity* $m$ *that these two components can contribute at (1,1) is given by*

$$m \leq (\lfloor \alpha_1 \rfloor + 1) + (\lfloor \alpha_2 \rfloor + \min\{\lfloor \frac{K+1}{2} \rfloor + \lfloor \delta_1 + \delta_2 \rfloor, \lfloor \alpha_2 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\})$$

$\square$

Theorem 7.5 implies that there is actually a combination of two components which yields a multiplicity between that which would be indicated by i) both components of non-binary type and ii) one component optimal binary and the other non-binary. In the same way double-binary components can also be suboptimal with a multiplicity greater than that indicated by Theorem 6.9 but less than that indicated by Theorem 6.15. Such components are again limited in multiplicity by the number of branch points they contain.

**Lemma 7.6** (Suboptimal Symmetric Double Binary, (cf. Lemma 6.13)). *Let* $\Omega_1$ *be a symmetric double-binary component containing both zero points,* $(0, w_1^0)$ *and* $(0, w_2^0)$, *with leading exponents* $-\alpha_1$ *and* $-\alpha_2$ *of* $\varphi$ *at* $(0, w_1^0)$ *and* $(0, w_2^0)$, *respectively. Further, suppose that* $\Omega_1$ *contains at most* $K$ *branch points of* $w(z, \mu)$. *Then the multiplicity* $m$ *of* $\Omega_1$ *satisfies*

$$m \leq \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor + \min\{\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1 + \lfloor \delta_1 + \delta_2 \rfloor, \lfloor \frac{K+1}{2} \rfloor\}.$$
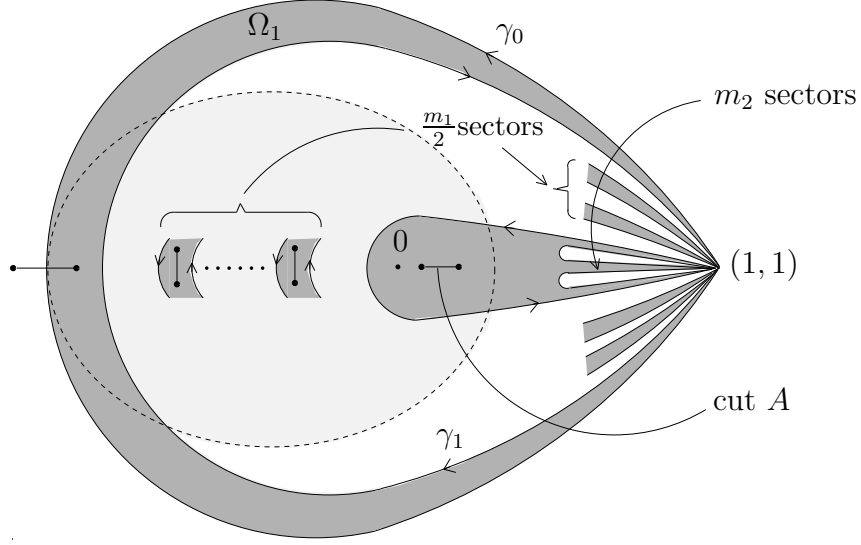
Figure 16: Suboptimal symmetric double binary component

Proof: The argument principle is applied assuming $m_2$ non-binary loops of $\Omega_1$ at $z = 1$, and $\frac{m_1}{2}$ binary loops of $\Omega_1$ at $z = 1$. In this case

$$m = m_2 + m_1 + 3,$$

and by the argument principle,

$$0 \leq \lfloor \alpha_1 + \alpha_2 \rfloor + \lfloor -\alpha_1 \rfloor + \lfloor \alpha_1 \rfloor - m_2 - \frac{m_1}{2}.$$

Therefore

$$m \leq \lfloor \alpha_1 + \alpha_2 \rfloor - 1 + 3 + \frac{m_1}{2},$$

where $m_1$ is limited by the total number of branch points $K$,

$$m_1 \leq K - 3.$$

Therefore,

$$m \leq \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor + \lfloor \frac{K+1}{2} \rfloor,$$

and the result follows in combination with Lemma 6.13. $\qquad \square$

Again there is a minimum number of branch points for which an optimal double-binary component can be obtained.

**Corollary 7.7.** *The minimum number of branch points, $K_{\min}$, contained in a component $\Omega_1$, described by Lemma 7.6, for which the maximum multiplicity, as indicated by Lemma 6.13, is obtained, is given by*

$$\lfloor \frac{K_{\min}+1}{2} \rfloor = \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + \lfloor \delta_1 + \delta_2 \rfloor + 1.$$

$\square$

27

Figure 17: Suboptimal non-symmetric double binary component

The non-symmetric double binary component, as illustrated by Fig. 13, may also be suboptimal, see Fig. 17.

**Lemma 7.8** (Suboptimal Non-symmetric Double Binary, (cf. Lemma 6.14)). *Let $\Omega_1$ be a non-symmetric double-binary component containing both zero points, $(0, w_1^0)$ and $(0, w_2^0)$, with leading exponents $-\alpha_1$ and $-\alpha_2$ of $\varphi$ at $(0, w_1^0)$ and $(0, w_2^0)$, respectively. Further, suppose that $\Omega_1$ contains at most $K$ branch points of $w(z, \mu)$. Then the multiplicity $m$ of $\Omega_1$ satisfies*

$$m \leq \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1 + \min\{\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1, \lfloor \frac{K+1}{2} \rfloor\}.$$

**Corollary 7.9.** *The minimum number of branch points, $K_{\min}$, contained in a component $\Omega_1$, described as in Lemma 7.8, for which the maximum multiplicity, as indicated by Lemma 6.14, is obtained, is given by*

$$\lfloor \frac{K_{\min} + 1}{2} \rfloor = \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1 .$$

$\square$

We note that other geometries, for example as illustrated in Fig. 27, in the appendix A, may lead to less efficient use of the branch points. Combining the conclusions of Lemmas 7.6 and 7.8 we deduce the following theorem.

28

**Theorem 7.10** (Suboptimal Double Binary). *Let $\varphi$ have leading exponents of $-\alpha_1$ and $-\alpha_2$ at the zero points $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$, and suppose that both zero points belong to a double-binary component $\Omega_1$ which also contains $K$ branch points. Then the highest total multiplicity, $m$, that this component can contribute at (1,1) is given by*

$$m \leq \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + \max\{\min\{\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 2, \lfloor \frac{K+1}{2} \rfloor + 1\},$$

$$\min\{\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1 + 2\lfloor \delta_1 + \delta_2 \rfloor, \lfloor \frac{K+1}{2} \rfloor + \lfloor \delta_1 + \delta_2 \rfloor\}\},$$

$\square$

## 7.1 Proof of Conjecture 6.11 for Maximum Order Schemes

Here we derive bounds on the order, $p$, of explicit schemes, under the assumption that the double-binary components described in Sections 6.1 and 7, exist. But we note that their multiplicities depend very intimately on the number of branch points contained inside the unit disk. These components cannot use branch cuts completely outside the unit circle since then stability would be violated. We also know, by Lemma 5.2, that the total number, $m$, of sectors at $z = 1$, both from inside and from outside the unit disk, satisfies $p + 1 = m$. Let $m_I$ and $m_O$ be the number of sectors of $\Omega$ at $z = 1$ from inside, and outside the unit disk, respectively. Then $m = m_I + m_O$, and by Lemma 5.2,

$$m_O - 1 \leq m_I \leq m_O + 1.$$

Therefore, if considered independently

(18) $$p \leq \max\{2m_I, 2m_O\}.$$

But, since $m_O$ can be seen to be limited by the number of branch points of $w(z, \mu)$ outside the unit disk, the bound on $m_O$ actually depends on the bound on $m_I$, via the number of branch points, $K_I$, inside the unit disk, limiting the number of branch points, $K_O$, outside the unit disk, because for convex schemes

$$2(r_1 + s_1) = K_I + K_O.$$

We will show that for the schemes of maximal order, $p = |I| - 2$, an argument in which both $m_I$ and $m_O$ are determined is required to derive a tight bound on $p$. But first we have to consider how to obtain $m_O$.

Instead of repeating the analyses of the earlier sections we apply a symmetry argument for the $\Omega$-sectors outside $\Delta$. For generality, we consider here implicit schemes. The functions $a_i$ defined in (6) will here be written in the form $a_i(z, \mu)$ to emphasize the $\mu$-dependence of the coefficients $a_{ij}$. Suppose we have a stable scheme (4) and that its stencil is regular for a certain value of $\mu$, with $R$ downwind and $S$ upwind stencil points
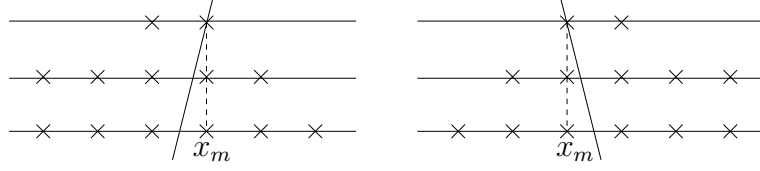
Figure 18: Stencil of reversed scheme

according to (9). Let $\ell$ denote the number of $\Omega$-sectors of the corresponding order star $\Omega$ emerging from the point (1,1) outside the unit disk $\Delta$. Then consider the reversed scheme

(19)
$$\sum_{j=-s_2}^{r_2} a_{2,-j}(-\mu)\, u_{n+2,m+j} + \sum_{j=-s_1}^{r_1} a_{1,-j}(-\mu) u_{n+1,m+j} +$$
$$\sum_{j=-s_0}^{r_0} a_{0,-j}(-\mu) u_{n,m+j} = 0 .$$

This scheme can be thought of as being obtained by a transformation of the space variable $x$ into $-x$. Hence the stencil is reflected about the line $x = x_m$.

The new scheme has $R^* = S$ downwind and $S^* = R$ upwind stencil points with respect to the characteristic $\mu^* = -\mu$. The characteristic function $\Phi^*$ of the reversed scheme is obtained from the characteristic function $\Phi$ by using the transformation

$$z \to \frac{1}{z} \text{ and } \mu \to -\mu .$$

Hence
$$\begin{aligned}
\Phi^*(z, w^*, \mu^*) &= \Phi(\tfrac{1}{z}, w^*, -\mu) \\
&= a_2(\tfrac{1}{z}, -\mu)(w^*)^2 + a_1(\tfrac{1}{z}, -\mu)w^* + a_0(\tfrac{1}{z}, -\mu) . \\
&= a_2^*(z, \mu^*)w^{*2} + a_1^*(z, \mu^*)w^* + a_0^*(z, \mu^*).
\end{aligned}$$

Therefore the algebraic function $w^*$ satisfies

(20)
$$w^*(z, \mu) = w(\tfrac{1}{z}, -\mu) .$$

From this relationship it follows that the reversed scheme is stable and of order $p$ if and only if the original scheme is stable and of order $p$. The order star $\Omega^*$ of the reversed scheme (19) is related to $\Omega$ of (4) in the sense that the portion of $\Omega$ outside the unit disk $\Delta$ is mapped to the inside of $\Delta$ and vice versa by the mapping $z \to \frac{1}{z}$.

Therefore, in order to determine $m_O$ for a given value of $\mu$, we should map $z \to \frac{1}{z}$ and investigate the portion of $\Omega$ inside $\Delta$ for $\mu^* = -\mu$. For explicit schemes the corresponding leading exponents of $\varphi$ will be $-\beta_1$ and $-\beta_2$ at $(0, w_1^0)$ and $(0, w_2^0)$, respectively, where $\beta_1 = s_0 - s_1 + \mu^*$, $\beta_2 = s_1 - s_2 + \mu^*$. Effectively, this maps the pole at infinity to zero, so that the effects of the infinity points can be examined via the effects of the zero points

for $\mu^*$. Note here that the form of $\beta_1$ and $\beta_2$ explicitly assumes $s_1 - s_2 \geq s_0 - s_1$, see Proposition 4.4 and Remark 4.5. Hence the argument we adopt enforces convexity on both sides of the characteristic line.

To complete the methodology we also need to be able to count the number of branch points of $w^*(z, \mu^*)$ inside the unit disk for $\mu^*$. But, by (20) this is just the number of branch points of $\omega(\frac{1}{z}, -\mu^*)$, which we have already denoted by $K_O$.

**Lemma 7.11** (Conjecture 6.11 for explicit maximal order schemes ). *Suppose we have an explicit stable scheme of type (4) of maximal order, $p = |I| - 2$, with a convex increasing stencil, with a fixed Courant number $\mu$, $-\frac{1}{2} < \mu < 0$. Assume that the algebraic function $w$ of $\Phi(z, w) = 0$ has no branch point at $z = 0$ and that $\varphi$ has leading exponents of $-\alpha_1$ and $-\alpha_2$ at the zero points $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$. Then the multiplicity $m$ of the $\Omega-$component $\Omega_1$ containing both zero points satisfies*

$$m \leq \lfloor \alpha_1 \rfloor + 1 + 2\lfloor \alpha_2 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\}.$$

**Proof:** Note that this is a statement of Conjecture 6.11 for the schemes of maximal order, $p_{opt} = |I| - 2$. Hence what we seek to prove is that for these schemes there cannot exist components of double-binary or suboptimal binary type. In particular we show that if these components exist the order is necessarily less that $p_{opt}$.

From Proposition 4.4 and Remark 4.5 we have the following expansions of $\varphi(z, w_i(z))$ at $z = 0$:

$$\begin{aligned}
\varphi(z, w_1(z)) &= z^{-r_1-\mu}(b_0 + b_1 z + b_2 z^2 + \ldots), \\
\varphi(z, w_2(z)) &= z^{-(r_0-r_1)-\mu}(c_0 + c_1 z + c_2 z^2 + \ldots).
\end{aligned}$$

It should be noted here that we have no means of associating a certain expansion with the zero point on a specific sheet. Hence the expansions will be associated with the zero points in the way which leads to the highest possible multiplicity.

Equivalently, for $\mu^* = -\mu$ we have the expansions of $\varphi(z, w_i^*(z))$ at $z = 0$

$$\begin{aligned}
\varphi(z, w_1^*(z)) &= z^{-s_1-\mu^*}(b_0^* + b_1^* z + b_2^* z^2 + \ldots), \\
\varphi(z, w_2^*(z)) &= z^{-(s_0-s_1)-\mu^*}(c_0^* + c_1^* z + c_2^* z^* + \ldots),
\end{aligned}$$

obtained via the transformation $z = \frac{1}{z}$. To avoid confusion we denote the exponents associated with $\mu$ by

$$\begin{aligned}
\alpha_1 &= r_0 - r_1 + \mu \\
\alpha_2 &= r_1 + \mu,
\end{aligned}$$

(21)

and those associated with $\mu^*$ by

$$\begin{aligned}
\beta_1 &= s_0 - s_1 + \mu^*, \\
\beta_2 &= s_1 + \mu^*.
\end{aligned}$$

(22)

31

By convexity, $\lfloor \alpha_2 \rfloor \geq \lfloor \alpha_1 \rfloor$ and $\lfloor \beta_2 \rfloor \geq \lfloor \beta_1 \rfloor$. Furthermore, we explicitly assume $\lfloor \alpha_1 \rfloor > 0$. Otherwise the zero point at $(0, w_1^0)$ does not lie inside $\Omega$ and the argument is considerably simplified. Similarly, assume $\lfloor \alpha_2 \rfloor > \lfloor \alpha_1 \rfloor$.

Clearly $-\frac{1}{2} < \mu < 0$ implies $\lfloor \delta_1 + \delta_2 \rfloor = \lfloor 2\delta_1 \rfloor = \lfloor 2(1+\mu) \rfloor = 1$.

(i) By Theorems 7.5 and 7.10 we see that the optimal configuration is dependent on the number of branch points utilised by the components. In particular define $K_L$ and $K_U$ to be the minimum number of branch points for which an optimal binary-non-binary configuration, and a double binary configuration is possible, respectively. Then

(23)
$$m_I \leq \begin{cases} \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 2 + \lfloor \frac{K_I+1}{2} \rfloor, & K_I < K_L \\ \lfloor \alpha_1 \rfloor + 1 + 2\lfloor \alpha_2 \rfloor + 2, & K_I = K_L \\ \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1 + \lfloor \frac{K_I+1}{2} \rfloor, & K_L < K_I < K_U \\ 2(\lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 1) + 1, & K_I \geq K_U. \end{cases}$$

Note that the bounds in (23) imply that the binary components are symmetric. Hence by corollaries 7.2 and 7.7

(24)
$$\lfloor \frac{K_L+1}{2} \rfloor = \lfloor \alpha_2 \rfloor + 1 \text{ and } \lfloor \frac{K_U+1}{2} \rfloor = \lfloor \alpha_1 \rfloor + \lfloor \alpha_2 \rfloor + 2.$$

Substituting for $\lfloor \alpha_1 \rfloor$ and $\lfloor \alpha_2 \rfloor$ in (23) and (24) we obtain

(25)
$$m_I \leq \begin{cases} r_0 + \lfloor \frac{K_I+1}{2} \rfloor, & \lfloor \frac{K_I+1}{2} \rfloor < r_1 \\ r_0 + r_1, & \lfloor \frac{K_I+1}{2} \rfloor = r_1 \\ r_0 - 1 + \lfloor \frac{K_I+1}{2} \rfloor, & r_1 < \lfloor \frac{K_I+1}{2} \rfloor < r_0 \\ 2r_0 - 1, & \lfloor \frac{K_I+1}{2} \rfloor \geq r_0. \end{cases}$$

(ii) Now we consider the outside of the unit disk and apply an equivalent argument with respect to $\beta_1, \beta_2$ and $\mu^* = -\mu$. In this case in the bound for $m_O$ we use $\lfloor \delta_1 + \delta_2 \rfloor = \lfloor 2\delta_1 \rfloor = \lfloor 2\mu \rfloor = 0$. Thus by Theorems 7.5 and 7.10

(26)
$$m_O \leq \begin{cases} \lfloor \beta_1 \rfloor + \lfloor \beta_2 \rfloor + 1 + \lfloor \frac{K_O+1}{2} \rfloor, & K_O < K_L \\ \lfloor \beta_1 \rfloor + 1 + 2\lfloor \beta_2 \rfloor + 1, & K_O = K_L \\ \lfloor \beta_1 \rfloor + \lfloor \beta_2 \rfloor + \lfloor \frac{K_O+1}{2} \rfloor + 1, & K_L < K_O < K_U \\ 2(\lfloor \beta_1 \rfloor + \lfloor \beta_2 \rfloor + 1), & K_O \geq K_U. \end{cases}$$

These components are non-symmetric and therefore, by corollaries 7.4 and 7.9,

(27)
$$\lfloor \frac{K_L+1}{2} \rfloor = \lfloor \beta_2 \rfloor + 1 \text{ and } \lfloor \frac{K_U+1}{2} \rfloor = \lfloor \beta_1 \rfloor + \lfloor \beta_2 \rfloor + 1.$$

As in (i), substitution of values for $\lfloor \beta_1 \rfloor$ and $\lfloor \beta_2 \rfloor$ in (26) and (27) leads to

(28)
$$m_O \leq \begin{cases} s_0 + 1 + \lfloor \frac{K_O+1}{2} \rfloor, & \lfloor \frac{K_O+1}{2} \rfloor < s_1 + 1 \\ s_0 + s_1 + 2, & \lfloor \frac{K_O+1}{2} \rfloor = s_1 + 1 \\ s_0 + \lfloor \frac{K_O+1}{2} \rfloor + 1, & s_1 + 1 < \lfloor \frac{K_O+1}{2} \rfloor < s_0 + 1 \\ 2(s_0 + 1), & \lfloor \frac{K_O+1}{2} \rfloor \geq s_0 + 1. \end{cases}$$

32

(iii) We now combine the results from inside and outside $\Delta$.

First, observe that if $K_I$ is even, so is $K_O$, and because the total number of branch points is $2(r_1 + s_1)$,

$$\lfloor \frac{K_I + 1}{2} \rfloor + \lfloor \frac{K_O + 1}{2} \rfloor = r_1 + s_1,$$

whereas, if $K_I$ and $K_O$ are odd,

$$\lfloor \frac{K_I + 1}{2} \rfloor + \lfloor \frac{K_O + 1}{2} \rfloor = r_1 + s_1 + 1 .$$

Therefore bounds on $K_I$ imply bounds on $K_O$, and vice versa. Hence, only certain combinations of components inside and outside $\Delta$ are possible.

In particular, suppose that inside $\Delta$ there is a double-binary or suboptimal double binary configuration. Then by (25) $\lfloor \frac{K_I+1}{2} \rfloor > r_1$ and

$$\lfloor \frac{K_O + 1}{2} \rfloor < s_1 + \begin{cases} 0 & K_I \text{ even} \\ 1 & K_I \text{ odd.} \end{cases}$$

Therefore outside $\Delta$ there can be at most a suboptimal binary configuration and

$$\begin{aligned} m_I + m_O \; &\leq r_0 - 1 + s_0 + 1 + \lfloor \frac{K_O+1}{2} \rfloor + \lfloor \frac{K_I+1}{2} \rfloor \\ &= r_0 + s_0 + r_1 + s_1 + \begin{cases} 0 & K_I \text{ even.} \\ 1 & K_I \text{ odd} \end{cases} \end{aligned}$$

Hence

$$m_I + m_O \leq |I| - 3 + \begin{cases} 0 & K_I \text{ even} \\ 1 & K_I \text{ odd.} \end{cases}$$

By Lemma 5.2, therefore

$$p \leq |I| - 3,$$

and (15) follows for $-\frac{1}{2} < \mu < 0$. $\qquad \square$

We could now repeat the arguments for $0 < \mu < \frac{1}{2}$ but it is sufficient to examine the order star outside $\Delta$ for $-\frac{1}{2} < \mu < 0$. Hence again we would like to show that the bounds on $m_O$, (28), for double binary components, lead to a contradiction. For these components $\lfloor \frac{K_O+1}{2} \rfloor \geq s_1 + 1$, and hence $\lfloor \frac{K_I+1}{2} \rfloor \leq r_1$. Thus by (25)

$$m_I \leq r_0 + \lfloor \frac{K_I + 1}{2} \rfloor,$$

and

$$\begin{aligned} m_O + m_I \; &\leq r_0 + s_0 + \lfloor \frac{K_I+1}{2} \rfloor + \lfloor \frac{K_O+1}{2} \rfloor + 1 \\ &\leq r_0 + s_0 + r_1 + s_1 + 1 + \begin{cases} 0 & K_O \text{ even} \\ 1 & K_O \text{ odd.} \end{cases} \end{aligned}$$

This time

$$m_O + m_I \leq |I| - 2 + \begin{cases} 0 & K_O \text{ even} \\ 1 & K_O \text{ odd,} \end{cases}$$

33

and we do not obtain the required contradiction, unless it can be demonstrated that $K_O$ is even.

Note that at no point did we explicitly impose stability in the above proof, because $K_I$ and $K_O$ are simply the number of branch points used by $\Omega$ from inside and outside $\Delta$, respectively. But if a component inside $\Delta$ utilised a cut outside $\Delta$ then this would in fact require that the stability condition, Lemma 5.1, is violated. Hence $K_I$ and $K_O$ do actually refer to the number of branch points inside and outside $\Delta$, respectively, and stability is required.

# 8 Implicit schemes: the role of poles away from $z = 0$ on multiplicity

## 8.1 Components containing zero points and poles

Assume we have inside a binary component poles of a total multiplicity $p$. When applying the argument principle $-p$ is added on the left side of the equation, e.g. (11) or (12). This leads to a bound for $m$ with an additional term $2p$. If the component is nonbinary then this additional term is clearly only $p$. Hence, having poles in a nonbinary component is less efficient than having one in a binary component. However, we shall show that we then can get a contribution $3p$ if the poles have multiplicity 1. To be able to do this the component is not allowed to contain a zeropoint. Such components are treated in the next section.

## 8.2 Components containing only poles away from $z = 0$

We consider the relationship between the number of poles inside a component $\Omega_1$ and its multiplicity under the assumption that $\Omega_1$ contains neither zero point. Suppose a single pole $P$ of multiplicity $p$ lies inside $\Omega_1$.

**Example 8.1.** *The pole $P$ occurs on the positive real axis (for simplicity we locate it on the principal sheet in Fig. 19), or at any other location away from the real axis inside $\Delta$. In this case the zero points do not have any effect on the component $\Omega_1$, since the cuts $L_i$ can be chosen such that they do not interact with $\Omega_1$ in any way. Then $\Omega_1$ corresponds to the type of component treated in [24], where it was shown that the multiplicity m of $\Omega_1$ is bounded by*
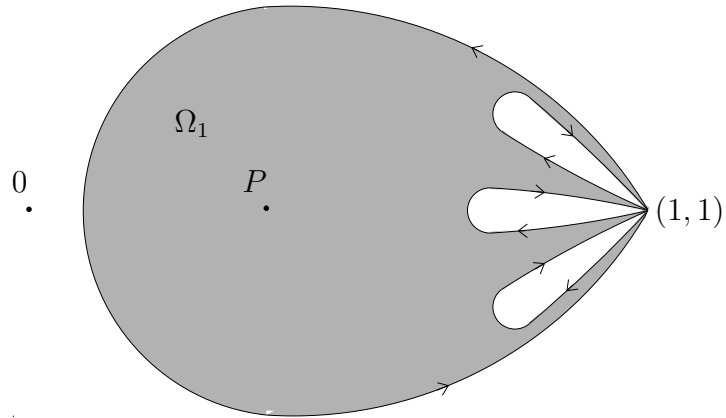
$$m \leq p \, .$$

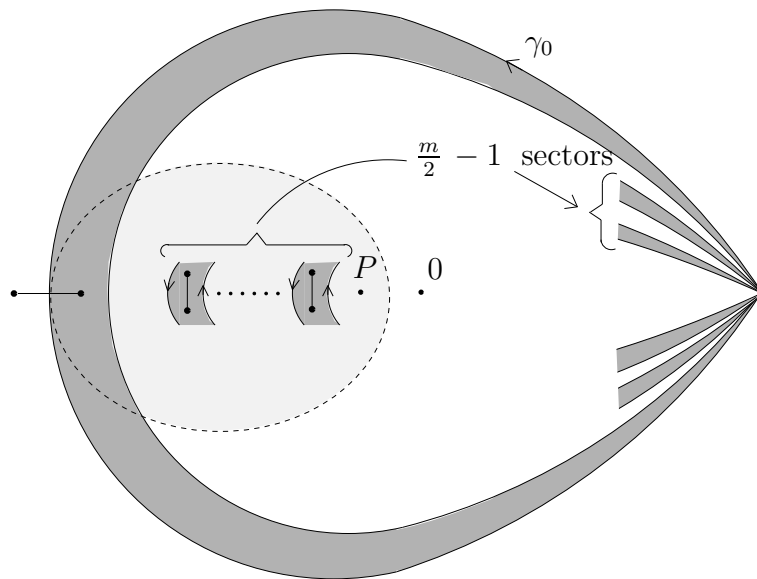Figure 19: Component with pole $P$ on positive real axis on principal sheet of $M$



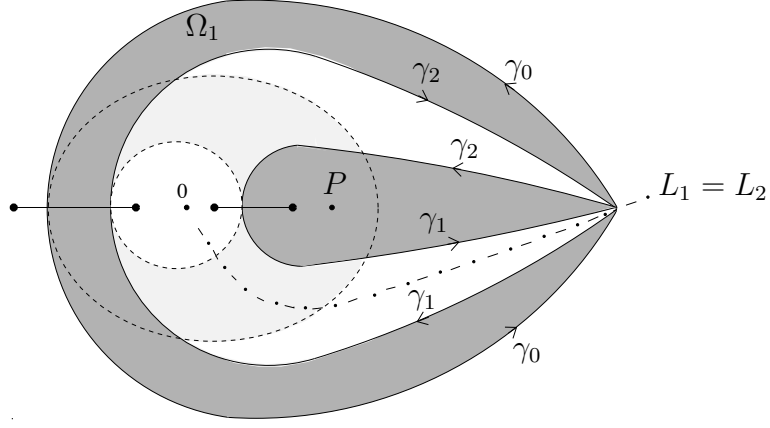Figure 20: Binary component with pole $P$ on negative real axis and one zeropoint encircled

Figure 21: Component with pole $P$ and with both zero points encircled

**Example 8.2.**  *For the components in Figs. 3, 4 the argument principle yields*

$$-p = \frac{1}{2\pi i} \int_{\gamma_0^E} + \frac{1}{2\pi i} \int_{\gamma_1^E} + \frac{1}{2\pi i} \sum_{j=1}^{m/2-1} \int_{\gamma_{2j}^E + \gamma_{2j+1}}$$

$$\leq \lfloor \alpha_1 \rfloor + \lfloor -\alpha_1 \rfloor - (m/2 - 1) .$$

*Hence, $m \leq 2p$ as expected.*

**Example 8.3.** *Let $\Omega_1$ be such as in Fig. 21, with the pole $P$ on the positive real axis on the secondary sheet of $M$ and both zero points excluded from $\Omega_1$. The branch cuts $L_i$ are chosen such that their projection onto the $z$-plane passes through $z = 1$. From one side of $L_1$ we obtain*

$$\frac{1}{2\pi i} \int_{\gamma_0^E + \gamma_1} \frac{\varphi'}{\varphi} \, dz \leq \lfloor \alpha_2 \rfloor, \; \frac{1}{2\pi i} \int_{\gamma_2^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor -\alpha_2 \rfloor$$

*and*

$$-p \leq \lfloor \alpha_2 \rfloor + \lfloor -\alpha_2 \rfloor - (m_1 - 2) .$$

*From the other side of $L_1$ we obtain*

$$\frac{1}{2\pi i} \int_{\gamma_2 + \gamma_1^E + \gamma_2^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor -\alpha_2 + \alpha_2 - 1 \rfloor .$$

*By applying the argument principle with the pole $P$ inside the component and $m_2$ sectors at $z = 1$, we obtain*

$$-p \leq -1 - (m_2 - 1) .$$

*A combination of these results leads to the following bound on the total multiplicity $m = m_1 + m_2$ of the component:*

36

Figure 22: Component with pole $P$ and with both zero points encircled

$$(29) \qquad\qquad m \leq 2p + 1 \ .$$

**Example 8.4.** *Let $\Omega_1$ be such as in Fig. 22, with the pole $P$ again occurring on the positive real axis on the secondary sheet of $M$ and both zero points excluded from $\Omega_1$. Then by choosing the branch cuts $L_i$ to go through $z = 1$ and by applying the argument principle in exactly the same way as in Example 8.3, we again obtain the bound $m \leq 2p+1$ occurring in (29).*

The question arises whether a pole $P$ away from $z = 0$ can yield multiplicity higher than in Example 8.3 and 8.4.

**Lemma 8.5** (Multiplicity of a single pole)**.** *Let $P$ be a pole of order $p$ away from $z = 0$ inside an $\Omega$-component $\Omega_1$ from which the two zero points are excluded by positively oriented portions of $\partial\Omega$ which encircle both zero points. Then the multiplicity $m$ of $\Omega_1$ is bounded by*

$$m \leq 2\,p + 1 \ .$$

$\square$

We can deduce from Lemma 8.5 that the highest possible multiplicity of a component, $\Omega_1$ relative to the order $p$ of a pole $P$ away from $z = 0$ inside it, is obtained if $p = 1$. Or, directly formulated: the most efficient poles away from $z = 0$ are simple poles. For such a simple pole we obtain the bound $m \leq 3$ on the multiplicity of the corresponding component. This entails that, if we have a normalized scheme which introduces into the corresponding order star poles of total order $p > 1$ inside $\Delta$, then the highest possible contribution of these poles to the number of $\Omega$-sectors inside $\Delta$ is obtained if these poles are simple and occur on the real axis. But then the complication occurs that the symmetry of components with respect to the real axis does not allow the simultaneous occurrence of
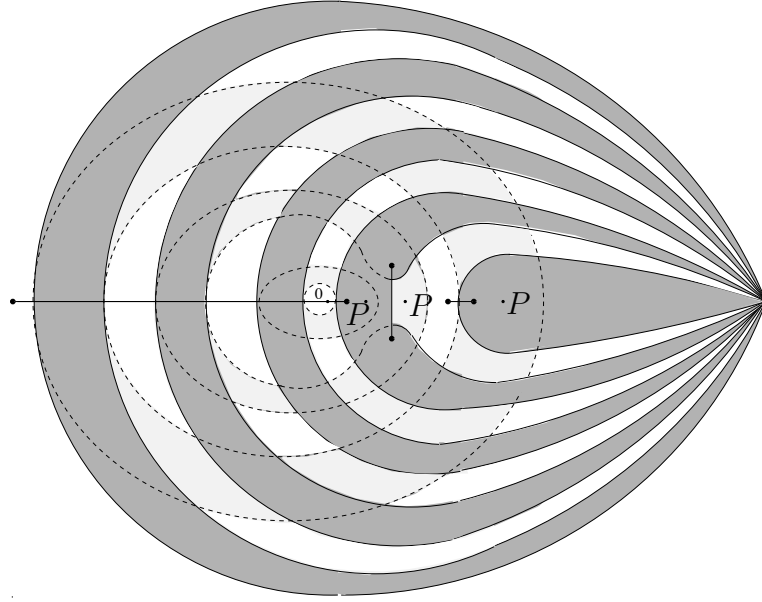
37

Figure 23: Three simple poles, each leading to multiplicity 3, inside $\Omega$-components

two separate components, each with a multiplicity of 3. As before this problem is overcome by means of two components of multiplicity 3 which are joined via a branch cut to yield one component of multiplicity 6. In this new component each simple pole still contributes 3 to the multiplicity of the component. This is illustrated in Fig. 23, where there are three simple poles, each leading to a contribution of 3 sectors to the total multiplicity inside $\Delta$. The rightmost pole belongs to a separate component, while the other two have joined to form a component of multiplicity 6. This situation is generalized.

**Proposition 8.6** (Maximum multiplicity of poles). *Let the order star of a stable, normalized scheme have poles of total multiplicity $p$ away from $z = 0$ inside $s\Delta$. Then the highest possible contribution of these poles to the multiplicity $m$ of components inside $\Delta$ is obtained if the poles are simple and real, leading to a multiplicity bounded by*

$$m \leq 3p \; .$$

$\square$

**Remark 8.7.** *Instead of two components with simple poles being joined to form one component such as in Fig. 23, a component containing a simple pole can also be joined with a component containing a zero point, such as in Fig. 24. There the combined multiplicity is 4. (The zero point $(0, w_2^0)$ belongs to a separate binary component of multiplicity 2). It can be seen that the efficiency of the pole and of the zero point $(0, w_1^0)$ remain unchanged as if they occur in separate components.*

**Remark 8.8.** *Observe from Figs. 21- 24 that in order for the poles in these components to each contribute a maximum multiplicity, of 3, branch cuts are required. In particular,*
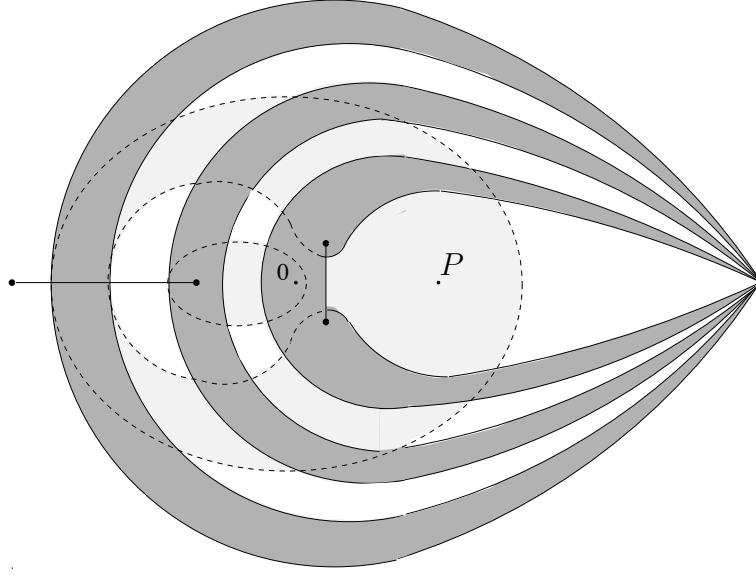
Figure 24: Component with a pole and one with a zero point which have been joined to form one component

*for a single pole a minimum of 3 branch points inside $\Delta$ are utilised. But for components containing more than one pole the branch points are used more efficiently. For two poles again just 3 branch points in $\Delta$ are sufficient. As further poles are added to the component each requires an additional cut inside $\Delta$. We therefore deduce the following corollary:*

**Corollary 8.9.** *Let the order star of a stable, normalized scheme have $n_p$ poles away from $z = 0$ inside a component $\Omega_1$, inside $\Delta$. Then the minimum number of branch points, $K_{\min}$, contained in $\Omega_1$ such that the multiplicity of $\Omega_1$ is given by $m = 3n_p$ satisfies*

$$K_{\min} = 2n_p - 1$$

$\square$

The question then arises as to whether Conjecture 6.11 can still be proved for the schemes of maximal order. This, however, turns out not to be so difficult. First let us consider both the order of the zero points and the number of branch points for the implicit schemes. When $r_2 > 0$, $s_2 > 0$ the exponents of $\varphi$ in the expansion around $z = 0$ are given by

$$\begin{aligned} \alpha_1 &= r_0 - r_1 + \mu \\ \alpha_2 &= r_1 - r_2 + \mu \end{aligned}$$

and

$$\begin{aligned} \beta_1 &= s_0 - s_1 + \mu^* \\ \beta_2 &= s_1 - s_2 + \mu^*, \end{aligned}$$

39

as compared with (21) and (22), respectively, when $r_2 = s_2 = 0$. But, by convexity, we still have $\lfloor \alpha_2 \rfloor \geq \lfloor \alpha_1 \rfloor$ and $\lfloor \beta_2 \rfloor \geq \lfloor \beta_1 \rfloor$, even though $\lfloor \alpha_2 \rfloor$ and $\lfloor \beta_2 \rfloor$ are reduced by $r_2$ and $s_2$, respectively. Hence the number of branch points $K$ utilised by the zero points is reduced by $2r_2 - 1$ if $K$ is odd and $2r_2$ if $K$ is even. But, by Corollary 8.9, the $r_2$ poles inside $\Delta$ need at least $2r_2 - 1$ branch points in order to contribute maximum multiplicity. Note further that for $K$ even, it was demonstrated in Lemma 7.11 that $p = |I| - 2$ could not be achieved. Hence the branch points left unutilised by the reduction of $\lfloor \alpha_2 \rfloor$ and $\lfloor \beta_2 \rfloor$, when $r_2, s_2 > 0$, are immediately required to contribute maximum multiplicity from the poles. Furthermore, since the contribution due to the poles gives a factor 3, rather than 2, in front of $n_p$ and $\lfloor \alpha_1 \rfloor$, respectively, we deduce that the optimal configuration uses branch points to maximise multiplicity due to the poles rather than due to the zero points. This leads us to conclude that Conjecture 6.11 is also valid for convex implicit schemes:

**Lemma 8.10** (Conjecture 6.11 for Implicit Schemes of Maximal Order).   *Suppose we have an implicit stable scheme of type (4) of maximal order, $p = |I| - 2$, with a convex increasing stencil, with a fixed Courant number $\mu$, $-\frac{1}{2} < \mu < 0$. Assume that the algebraic function $w$ of $\Phi(z, w) = 0$ has no branch point at $z = 0$ and that $\varphi$ has leading exponents of $-\alpha_1$ and $-\alpha_2$ at the zero points $(0, w_1^0)$ and $(0, w_2^0)$, respectively, on the two sheets of $M$. Then the multiplicity $m$ of the $\Omega-$component $\Omega_1$ containing both zero points and no poles satisfies*

$$m \leq \lfloor \alpha_1 \rfloor + 1 + 2\lfloor \alpha_2 \rfloor + \max\{1, 2\lfloor \delta_1 + \delta_2 \rfloor\}.$$

$\square$

# 9 Proof of the main theorem

The proof of Theorem 3.1 is divided into stages. One part of it is proved in Lemma 9.1 and the other part in Lemma 9.2. The numbers $R$ and $S$ denote the number of downwind and upwind stencil points, respectively, with respect to the characteristic through $(t_{n+2}, x_m)$ as defined in (9).

**Lemma 9.1** (Maximal order with stability). *Suppose we have a convex, normalized scheme of type (4) with a fixed Courant number $\mu$, $0 < |\mu| < \frac{1}{2}$. If the scheme is stable, and Conjecture 6.11 is satisfied, then the order $p$ of the scheme is bounded by*

$$(30) \qquad\qquad\qquad\qquad p \leq 2R \,.$$

**Proof:** Since the scheme is stable, there will be a clear distinction between the portion of the corresponding order star $\Omega$ inside $\Delta$ and the portion outside $\Delta$. In this proof we restrict ourselves to the portion inside $\Delta$.

From Proposition 4.4 we have the following expansions of $\varphi(z, w_i(z))$ at $z = 0$:

$$\varphi(z, w_1(z)) = z^{-(r_1-r_2)-\mu}(b_0 + b_1 z + b_2 z^2 + \ldots),$$
$$\varphi(z, w_2(z)) = z^{-(r_0-r_1)-\mu}(c_0 + c_1 z + c_2 z^2 + \ldots).$$

Again we have no means of associating a certain expansion with the zero point on a specific sheet. Hence the expansions will be associated with the zero points in the way which leads to the highest possible multiplicity.

In the remainder of the proof we have to work separately with the cases where $\mu < 0$ and where $\mu > 0$. Note also that where we assume Conjecture 6.11, Lemmas 7.11 and 8.10 give the result for $p = p_{opt}$, and $-\frac{1}{2} < \mu < 0$.

  a) We first assume $-\frac{1}{2} < \mu < 0$. Then the following choices of the indices $r_0, r_1, r_2$ lead to different combinations of $\Omega$-components inside $\Delta$.

    (i) $r_0 = r_1 = r_2 = 0$. Then also $R = 0$. According to the Courant-Friedrichs-Lewy condition the scheme cannot be convergent, i.e. it is impossible to have order $p \geq 1$ and stability simultaneously.

    (ii) $r_0 = r_1 = r_2 > 0$. There are $r_2$ poles away from $z = 0$ inside $\Delta$, while both $\varphi(z, w_1(z))$ and $\varphi(z, w_2(z))$ have positive leading exponents of $-\alpha_i = -\mu$ at $z = 0$, implying that both zero points belong to $\Omega^c$. We apply Proposition 8.6 to obtain

$$m \leq 3r_2 = r_0 + r_1 + r_2 = R.$$

    (iii) $0 = r_0 - r_1 < r_1 - r_2$. Then $-\alpha_1 = -\mu > 0$, implying that $(0, w_1^0) \in \Omega^c$, and $-\alpha_2 = -(r_1 - r_2) - \mu < 0$, implying that $(0, w_2^0) \in \Omega$. The highest possible multiplicity is obtained if $(0, w_2^0)$ belongs to a binary component, in which case we apply Proposition 6.6. If $r_2 > 0$, the poles away from $z = 0$ are again treated according to Proposition 8.6. Then we have

$$\begin{aligned} m &\leq 3r_2 + 2\lfloor r_1 - r_2 + \mu \rfloor + 2\lfloor 2 + 2\mu \rfloor \\ &= 3r_2 + 2(r_1 - r_2 - 1) + 2 \\ &= r_2 + 2r_1 = r_2 + r_1 + r_0 = R. \end{aligned}$$

    (iv) $0 < r_0 - r_1 < r_1 - r_2$. Then $-\alpha_1 = -(r_0 - r_1) - \mu < 0$ and $-\alpha_2 = -(r_1 - r_2) - \mu < 0$, implying that both $(0, w_1^0)$ and $(0, w_2^0)$ belong to $\Omega$. Since $\alpha_1 < \alpha_2$, the highest possible multiplicity is obtained by applying Conjecture 6.11, with $(0, w_1^0)$ inside a non-binary and $(0, w_2^0)$ inside a binary component. If $r_2 > 0$, the poles away from $z = 0$ are again treated according to Proposition 8.6. Then the total multiplicity $m$ inside $\Delta$ is bounded by

$$\begin{aligned} m &\leq 3r_2 + \{2\lfloor r_1 - r_2 + \mu \rfloor + 2\lfloor 2 + 2\mu \rfloor\} + \{\lfloor r_0 - r_1 + \mu \rfloor + 1\} \\ &= 3r_2 + \{2(r_1 - r_2 - 1) + 2\} + \{(r_0 - r_1 - 1) + 1\} \\ &= r_0 + r_1 + r_2 = R. \end{aligned}$$

41

b) Assume $0 < \mu < \frac{1}{2}$. Then we have $-\alpha_1 = -(r_0 - r_1) - \mu < 0$ and $-\alpha_2 = -(r_1 - r_2) - \mu < 0$, implying that both $(0, w_1^0)$ and $(0, w_2^0)$ belong to $\Omega$. Since $\alpha_1 \leq \alpha_2$, the highest possible multiplicity is obtained, assuming Conjecture 6.11, if $\alpha_1$ is inside a non-binary and $\alpha_2$ inside a binary component. If $r_2 > 0$, the poles away from $z = 0$ are treated according to Proposition 8.6. This leads to the bound

$$
\begin{aligned}
m &\leq 3r_2 &+& \{2\lfloor r_1 - r_2 + \mu \rfloor + 1\} + \{\lfloor r_0 - r_1 + \mu \rfloor + 1\} \\
&= 3r_2 &+& \{2(r_1 - r_2) + 1\} + \{(r_0 - r_1) + 1\} \\
&= r_0 &+& r_1 + r_2 + 2 = R \ .
\end{aligned}
$$

In all the foregoing cases we obtained

$$ m \leq R \ . $$

The remainder of the proof makes use of Lemma 5.2 and hence Equation (18) to give for the order $p$ of the scheme

$$ p + 1 \leq m + (m + 1) \leq 2R + 1 \ , $$

which leads to $p \leq 2R$. $\qquad \square$

Concerning the upwind points of a difference stencil we now prove the following lemma.

**Lemma 9.2** (Maximum order with stability). *Suppose we have a convex and normalized scheme of type (4) with a fixed Courant number $\mu$ satisfying $0 < |\mu| < \frac{1}{2}$. If the scheme is stable, then the order $p$ of the scheme is bounded by*

$$ p \leq 2S \ . $$

**Proof:** Instead of repeating the argument of Lemma 9.1 for the $\Omega$-sectors outside $\Delta$, we apply the symmetry argument introduced in Section 7 to prove this result.

Hence, if (30) is proved for a value of $\mu$ for which the stencil is regular, we obtain by the mapping $z \to \frac{1}{z}$ and $\mu \to -\mu$ for $-\mu$ that

$$ \ell \leq R^* = S \ . $$

By making use of Lemma 5.2 this result leads to

$$ p \leq 2S \ . $$
$\qquad \square$

**Acknowledgements:**

# References

[1] Courant R., Friedrichs K.O., Lewy H., Über die partiellen Differenzengleichungen der mathematischen Physik, Math. Ann. **100**, 32-74, (1928).

[2] Goldberg M., Simple stability criteria for difference approximations of hyperbolic initial-boundary value problems II, in: Proceedings of the Third International Conference on Hyperbolic Problems, Uppsala, June 11-15, 1990, ed. Engquist B. and Gustafsson B., Chartwell-Bratt, 519-527, (1991).

[3] Goldberg M., Tadmor E., Convenient stability criteria for difference approximations of hyperbolic initial-boundary value problems II, Math. Comp. **48**, 503-520, (1987).

[4] Gustafsson B., Kreiss H.-O., Sundström A., Stability theory of difference approximations for mixed initial boundary value problems II, Math. Comp. **26**, 649-686, (1972).

[5] Hairer E., Wanner G., Solving ordinary differential equations II, Springer-Verlag, 1991.

[6] Iserles A., Order and stability of fully-discretized hyperbolic finite differences, Report DAMPT 1985/NAG of the University of Cambridge, (1985).

[7] Iserles A., Order stars and a saturation theorem for first order hyperbolics, IMA J. Num. Anal. **2**, 49-61, (1982).

[8] Iserles A., Nørsett S. P., Order stars, Chapman & Hall, 1991.

[9] Iserles A., Strang G., The optimal accuracy of difference schemes, Trans. Amer. Math. Soc. **277**, 779-803, (1983).

[10] Jeltsch R., Stability and accuracy of difference schemes for hyperbolic problems, J. Comput. Appl. Math. **12 & 13**, 91-108, (1985).

[11] Jeltsch R., Order barriers for difference schemes for linear and non-linear hyperbolic problems, in: Numerical Analysis, Proceedings of the 1987 Dundee conference, ed. D. F. Griffiths, D. A. Watson, Pitman, (1988).

[12] Jeltsch R., Kiani P., Stability of a family of multi-time-level difference schemes for the advection equation, Numer. Math **60**, 77-95 (1991).

[13] Jeltsch R., Kiani P., Raczek K., Counterexamples to a stability barrier, Numer. Math. **52**, 301-316, (1988).

[14] Jeltsch R., Renaut R. A., Smit J. H., The maximal accuracy of stable difference schemes for the wave equation, Report of Forschungsinstitut für Mathematik, ETH Zürich, (1993), BIT 35, 1, 83-115 (1995).

[15] Jeltsch R., Smit J. H., Accuracy barriers of difference schemes for hyperbolic equations, SIAM J. Numer. Anal. **24**, 1-11, (1987).

[16] Jeltsch R., Smit J. H., Accuracy barriers of three-time-level difference schemes for hyperbolic equations, in: Proceedings of the Third International Conference on Hyperbolic Problems, Uppsala, June 11-15, 1990, ed. Engquist B. and Gustafsson B., Chartwell-Bratt, 611-627, (1991).

[17] Jeltsch R., Smit J. H., Accuracy barriers of three-time-level difference schemes for hyperbolic equations, Ann. University of Stellenbosch, 1992/2, 1-34, (1992).

[18] Kreiss H. O., Difference approximations for the initial-boundary value problem for hyperbolic differential equations, in: Numerical Solutions of Nonlinear Differential Equations. Proc. Adv. Sympos., Madison, Wisconsin, 141-166, (1966).

[19] Kreiss H. O., Stability theory for difference approximations of mixed initial-boundary value problems I, Math. Comput. **22**, 703-714, (1968).

[20] Renaut R. A., Full discretizations of $u_{tt} = u_{xx}$ and rational approximations to $\cos h\mu Z$, SIAM J. Numer. Anal. **26**, (1989).

[21] Renaut R. A., Smit J. H., Order stars and the maximal accuracy of stable difference schemes for the wave equation, Quaestiones Math. 15, 307-323, (1992).

[22] Smit J. H., Order stars and the optimal accuracy of stable, explicit difference schemes, Quaestiones Math. 8, 167-188, (1985).

[23] Strang G., Iserles A., Barriers to stability, SIAM J. Numer. Anal. **20**, 1251-1257, (1983).

[24] Wanner G., Hairer E., Nørsett S. P., Order stars and stability theorems, BIT 18, 475-489, (1978).

# Appendix A. Motivation of Conjecture 6.11

Here we return to the consideration of the multiplicities of double binary components. Observe that if cut $B$ in Fig. 13 would not be present and the curves $\gamma_0$ and $\gamma_1$ would join together as in Fig. 25 we would still obtain bound (19). Note that Fig. 25 corresponds to Fig. 26 c) of [16]. Now let us deform the boundary of $\Omega_1$ in Fig. 13. As long as the deformation is continuous this should not affect the argument principle if neither a zeropoint nor a zero crosses the boundary. Hence, the integrals along the curves given in Figs. 26-27 are the same since a branchpoint on a Riemann surface is not a special point with respect to integration. The connectivity of the curves in Fig. 26 b) can be interpreted in two ways, the one obtained from a continuous deformation of Fig. 26 a) and the one obtained from Fig. 27 by a continuous deformation. However, what does change in our way of applying the argument principle, when making this continuous deformation, is the pairing of the integrals. Hence, Fig. 26 a) leads to bound (17) while Fig. 27 leads to the bound

$$m \leq 2\lfloor \alpha_1 \rfloor + 2\lfloor \alpha_2 \rfloor + 2\lfloor \delta_1 + \delta_2 \rfloor \ .$$

Again the factor 2 in the term with $\lfloor \alpha_1 \rfloor$ is present but note now that this bound is sharper than that given by (17) and at the same time requires fewer branch cuts to obtain the multiplicity. However, the cut $A$ has now a similar function as cut $A$ in Fig. 12. If it would be removed $\Omega_1$ would be separated into two disconnected $\Omega$ components, one which is binary and one which is non-binary. This would again give the bound of Theorem 6.9. We can actually remove the cut $A$ in Figs. 12, 26 a), respectively, by deforming the boundary as indicated in Figs. 28.

In Fig. 28 b) when binary loops are introduced the set $\Omega$ has been decomposed into two separate components, a binary and a nonbinary giving again the bound of Theorem 6.9. In a similar fashion one can move the cut $A$ in Fig. 12 outside the order star and thus separate $\Omega_1$ into two disconnected components, see Figs. 29 - 31.

The problem with these continuous deformations is that one has to be able to show that they can be performed while keeping the other conditions of the schemes unchanged such as stability, instability, real coefficients, error order, difference stencil at least on the downwind side. While such a proof seems to be technically extremely difficult the possible deformations suggested above support Conjecture 6.11.

We do observe, however, that for all double binary components there is one cut which is required to maintain the connection between the components. We have called this cut, cut A. To be effective cut A has to allow, in the application of the argument principle, that the multiplicity due to the zero point on one sheet is "carried" via the cut to the other sheet and then back to the original sheet by cuts inside binary loops. If the cut is excluded by integrating around it, in the application of the argument principle, the components are disconnected. But it appears it is still not possible to draw any conclusions about whether the multiplicity can be "carried" in this way because the integral around the cut amounts to an integration around a "hole" of the surface, which is not uniquely defined. We do feel that a further consideration of this approach will lead to a proof of Conjecture

6.11.

## Appendix B. Components with one zero point which are basically different

The components we considered in the body of the paper have boundary curves which circle at most once around a zero point. The question arises as to what the effect will be if one or more boundary curves circle several times around the zero points.

Because a zero point has the highest possible efficiency if it occurs inside a binary component, we initially consider binary components. The first one is a symmetric binary component depicted in Fig. 32.

By choosing the branch cuts $L_i$ according to convention, the way in which the boundary curves cross the branch cuts can be represented schematically in Fig. 33.

Integration along the boundary curves leads to

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor 2\alpha_1 + 2\alpha_2 \rfloor \, ; \; \frac{1}{2\pi i} \int_{\gamma_1^E} \frac{\varphi'}{\varphi} \, dz \leq \lfloor -2\alpha_1 - \alpha_2 \rfloor \, .$$

If we apply the argument principle and take into account the binary character of the component such as we did in Lemma 6.3, we obtain

(31) $$0 \leq 2\{\lfloor 2\alpha_1 + 2\alpha_2 \rfloor + \lfloor -2\alpha_1 - \alpha_2 \rfloor\} - (m - 2) \, .$$

In order to simplify (31), let $\widetilde{\alpha} = 2\alpha_1 + \alpha_2$, $\widetilde{\delta} = \widetilde{\alpha} - \lfloor \widetilde{\alpha} \rfloor$ and $\delta_2 = \alpha_2 - \lfloor \alpha_2 \rfloor$ as before. Then we obtain

$$0 \leq \{\lfloor \widetilde{\alpha} + \alpha_2 \rfloor + \lfloor -\widetilde{\alpha} \rfloor\} - \frac{(m - 2)}{2} \, ,$$

leading to the bound

$$m \leq 2\lfloor \alpha_2 \rfloor + 2\lfloor \widetilde{\delta} + \delta_2 \rfloor$$

which is basically the same as in Lemma 6.3. The simplification in (31) is possible because the "outward" boundary curve $\gamma_0$ is "accompanied" by an "inner" boundary curve $\gamma_1$, ensuring that the component is well defined as it "switches" to and fro between the two sheets of $M$.

**Remark** *There are two variations to the component in Fig. 32 which deserve special attention.*

   a) *Suppose we have a component with an "outward" boundary $\gamma_0$ which corresponds to that in Fig. 32, while the "inner" boundary curve $\gamma_1$ does not cross one of the branch cuts on the real axis, such as AB in Fig. 34. Then at least one other "inner" boundary curve $\gamma_2$ (say) is needed for the component to be well defined. Since the curve $\gamma_2$ does not pass through (1,1) it is not contributing to the multiplicity of the*

*component. In [17] a curve of this kind is called an **inefficient curve** (see [17], Fig. 18). Integration along the boundary now leads to*

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'}{\varphi}\, dz \le \lfloor 2\alpha_1 + 2\alpha_2 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_1^E} \frac{\varphi'}{\varphi}\, dz \le \lfloor -\alpha_1 \rfloor,$$

$$\frac{1}{2\pi i} \int_{\gamma_2^E} \frac{\varphi'}{\varphi}\, dz \le \lfloor -\alpha_1 - \alpha_2 \rfloor \, .$$

*Because $\lfloor -\alpha_1 \rfloor + \lfloor -\alpha_1 - \alpha_2 \rfloor \le \lfloor -2\alpha_1 - \alpha_2 \rfloor$, the inequality involving the argument principle can be simplified to (31), leading to the same bound as before. Since in general $\lfloor \alpha \rfloor + \lfloor \beta \rfloor \le \lfloor \alpha + \beta \rfloor$, a simplification of this kind can be made wherever an inefficient curve "replaces" an efficient one. For this reason we will as far as possible make use of efficient curves.*

b) *Suppose we have a component with an "outward" boundary $\gamma_0$ which corresponds to that in Fig. 32, but with more than one efficient "inner" boundary curve such as $\gamma_1, \gamma_2$ and $\gamma_3$ in Fig. 35. Then a careful combination of the terms in the inequality involving the argument principle leads to an inequality similar to (31). (For this simplification it is convenient to make the assumption that the non-integer parts of $\alpha_1$ and $\alpha_2$ coincide, i.e. that $\delta_1 = \delta_2$ where $\delta_i = \alpha_i - \lfloor \alpha_i \rfloor$, $i = 1, 2$. From a practical point of view this assumption is not restrictive at all, since the factor $z^{-\mu}$ in $\varphi$ is throughout responsible for the non-integer parts of $\alpha_1$ and $\alpha_2$. However, since this assumption is not needed in the rest of the discussion and since a component of the type in Fig. 35 is most unlikely to occur in any "real" order star, we do our analysis without making this assumption).*

The situation in (31) can be generalized to the case where the two boundary curves "switch" several times from the one sheet to the other. It can again be represented schematically as in Fig. 36.

Integration along the boundary curves now leads to

$$\frac{1}{2\pi i} \int_{\gamma_0^E} \frac{\varphi'}{\varphi}\, dz \le \lfloor k\alpha_1 + k\alpha_2 \rfloor, \quad \frac{1}{2\pi i} \int_{\gamma_1^E} \frac{\varphi'}{\varphi}\, dz \le \lfloor -k\alpha_1 - (k-1)\alpha_2 \rfloor \, .$$

If we apply the argument principle and take into account the binary character of the component, we obtain

(32) $$0 \le \{ \lfloor k\alpha_1 + k\alpha_2 \rfloor + \lfloor -k\alpha_1 - (k-1)\alpha_2 \rfloor \} - \frac{(m-2)}{2} \, .$$

The result (32) is simplified by taking $\widetilde{\alpha} = k\alpha_1 + (k-1)\alpha_2$, $\widetilde{\delta} = \widetilde{\alpha} - \lfloor \widetilde{\alpha} \rfloor$. Then we obtain as in (31)

$$m \le 2\lfloor \alpha_2 \rfloor + 2\lfloor \widetilde{\delta} + \delta_2 \rfloor \, .$$

We finally consider the case of a non-symmetric binary component such as the one depicted in Fig. 37. The presence of a branch cut such as $AB$ now presents the complication that one more "inner" boundary curve $\gamma_3$ (say), is needed for the component to be well defined. With respect to the curve $\gamma_3$ we have the following possibilities:

i) $\gamma_3$ is inefficient and $\Omega_1$ is a non-symmetric binary component such as in Fig. 37;

ii) $\gamma_3$ passes through (1,1) to be efficient and $\Omega_1$ becomes a symmetric binary component.

In both cases it can be shown that the bound in Proposition 6.6 remains valid.

In the foregoing examples we considered components with more than one "outward" and more than one "inner" boundary curve, being symmetric and non-symmetric with respect to the real axis and circling once or more than once around the zero point. In all these cases the bound on the multiplicity coincides with the result in Proposition 6.6. A similar analysis with respect to non-binary components leaves Proposition 6.1 unchanged.

It should also be noted that binary components such as those in Figs. 32, 34, 35, 37 cannot be combined with a non-binary component with odd multiplicity such as in Fig. 11.

In view of these considerations we have only worked with components with boundary curves which circle once around one or two zero points.

Figure 25: Simplification of non-symmetric binary component



Figure 26: a, b Deformation of the boundary curves of non-symmetric binary component
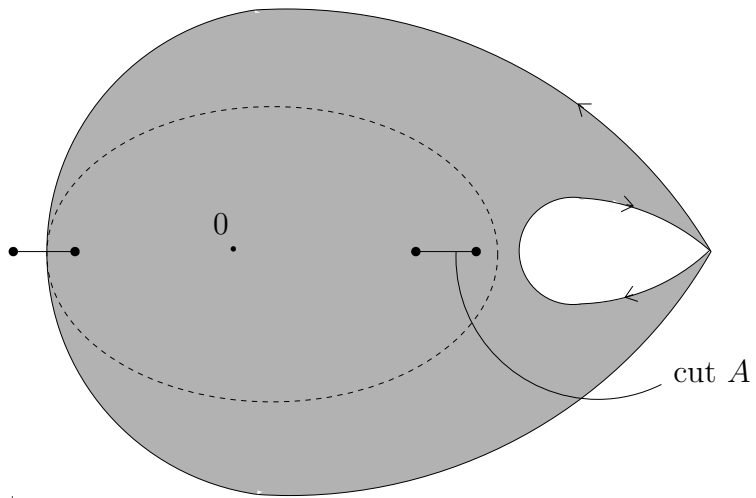


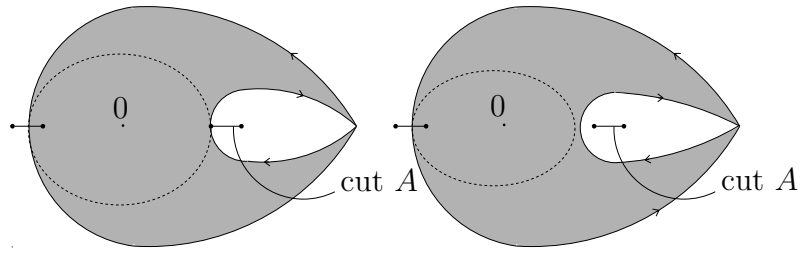Figure 27: Cut A moved inside the component

Figure 28: a, b Removal of cut A from the component
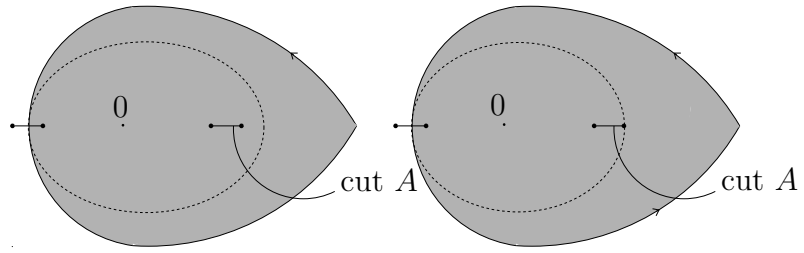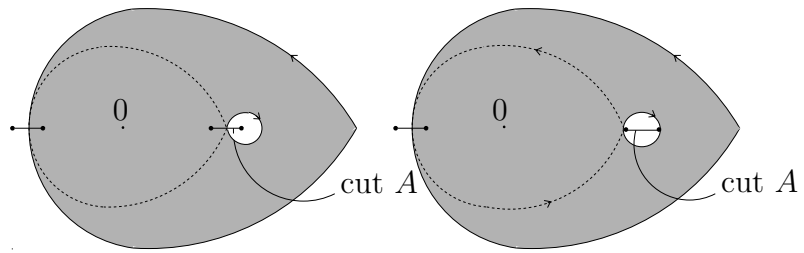


Figure 29: a, b Cut A redundant



Figure 30: Cut A being moved outside the order star

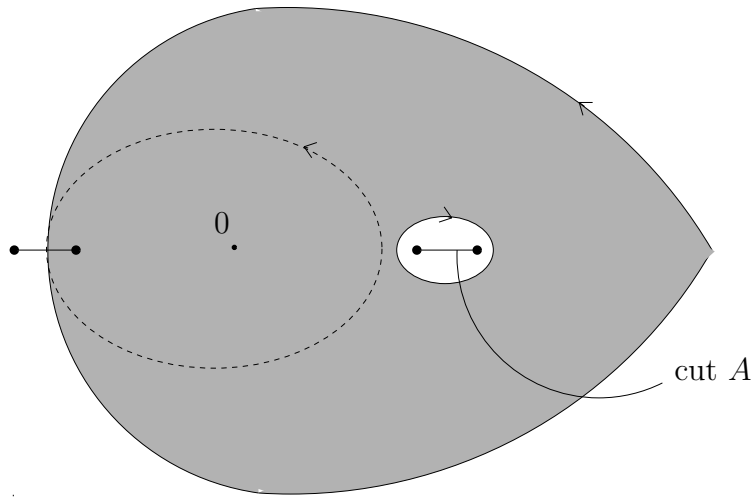Figure 31: Cut A outside the order star

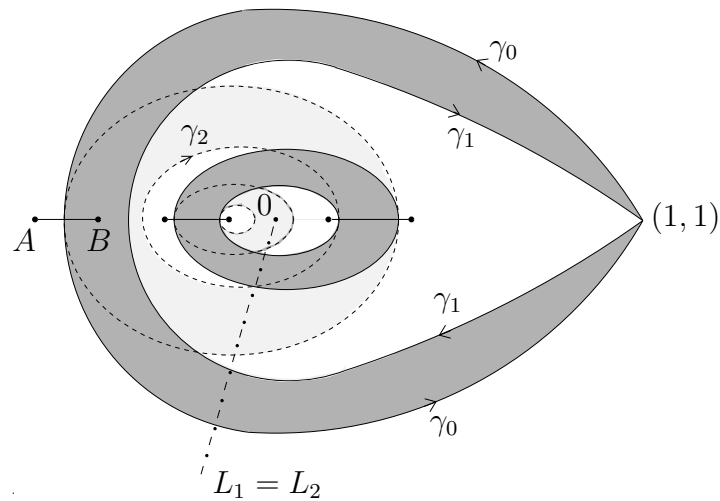Figs. 29 - 31    Transformation of $\gamma_0$ in Fig. 12 such as to move cut $A$ outside the order star.



Figure 32: Symmetric binary component with boundary curves circling more than once around the zero points
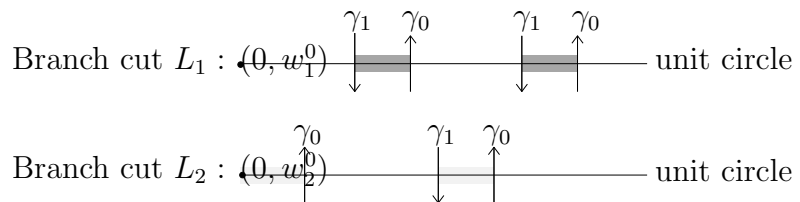


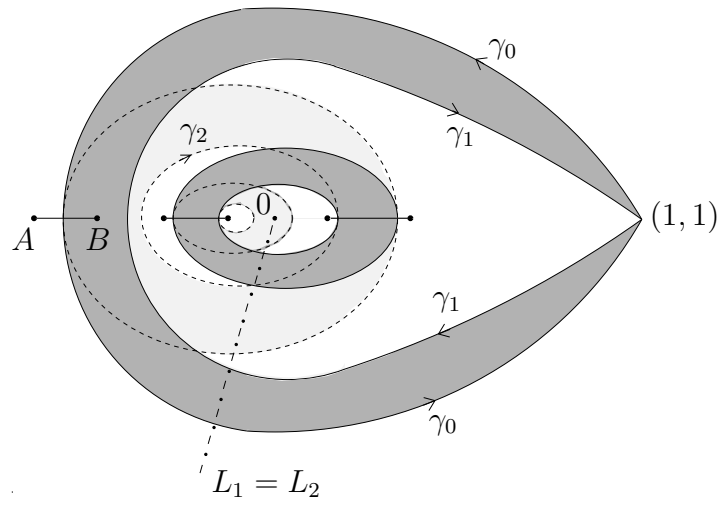Figure 33: Boundary curves crossing the branch cuts

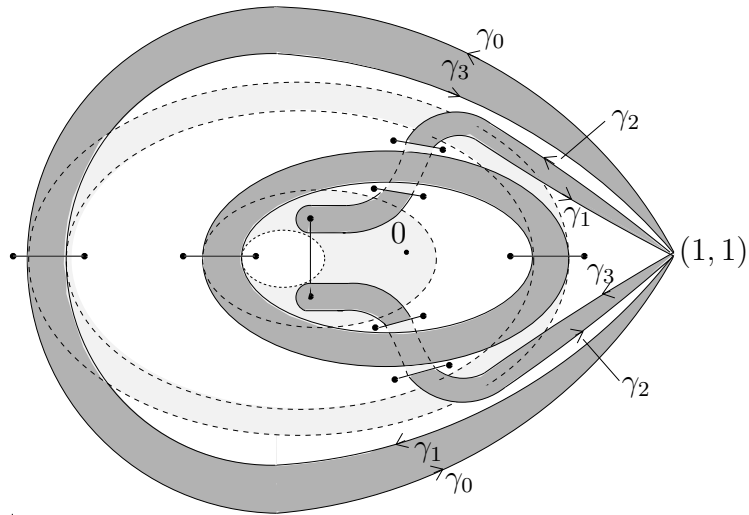Figure 34: Boundary curve $\gamma_2$ which is inefficient



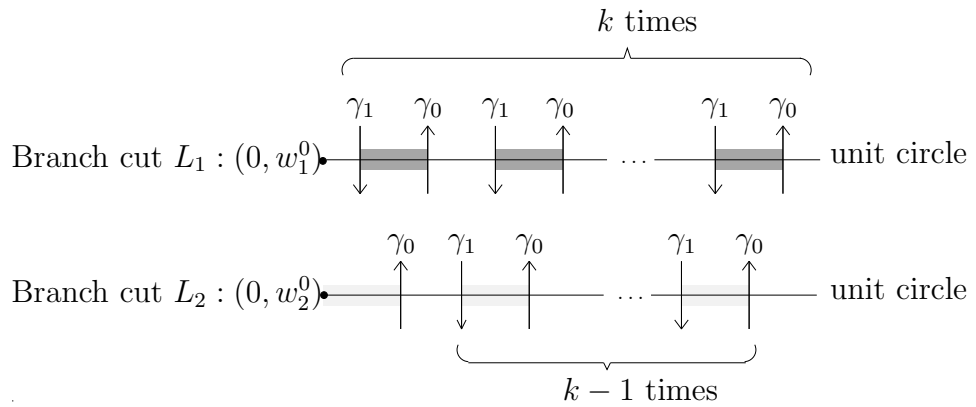Figure 35: More than one efficient "inner" boundary curve circling more than once around the zero points

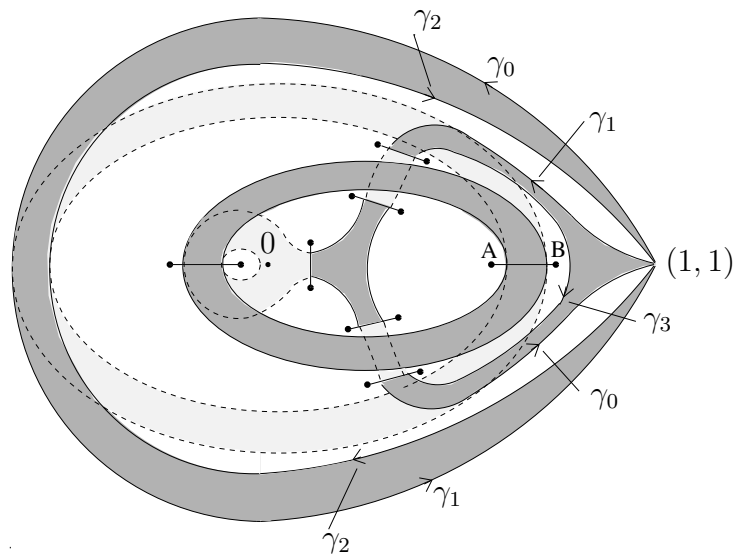Figure 36: Two boundary curves switching from one sheet to the other



Figure 37: Non-symmetric binary component with boundary curves circling more than once around the zero points

# Research Reports

| No. | Authors | Title |
| --- | --- | --- |
| 97-10 | R. Jeltsch, R.A. Renaut, J.H. Smit | An Accuracy Barrier for Stable Three-Time-Level Difference Schemes for Hyperbolic Equations |
| 97-09 | K. Gerdes, A.M. Matache, C. Schwab | Analysis of membrane locking in $hp$ FEM for a cylindrical shell |
| 97-08 | T. Gutzmer | Error Estimates for Reconstruction using Thin Plate Spline Interpolants |
| 97-07 | J.M. Melenk | Operator Adapted Spectral Element Methods. I. Harmonic and Generalized Harmonic Polynomials |
| 97-06 | C. Lage, C. Schwab | Two Notes on the Implementation of Wavelet Galerkin Boundary Element Methods |
| 97-05 | J.M. Melenk, C. Schwab | An $hp$ Finite Element Method for convection-diffusion problems |
| 97-04 | J.M. Melenk, C. Schwab | $hp$ FEM for Reaction-Diffusion Equations. II. Regularity Theory |
| 97-03 | J.M. Melenk, C. Schwab | $hp$ FEM for Reaction-Diffusion Equations. I: Robust Exponentiel Convergence |
| 97-02 | D. Schötzau, C. Schwab | Mixed $hp$-FEM on anisotropic meshes |
| 97-01 | R. Sperb | Extension of two inequalities of Payne |
| 96-22 | R. Bodenmann, A.-T. Morel | Stability analysis for the method of transport |
| 96-21 | K. Gerdes | Solution of the $3D$-Helmholtz equation in exterior domains of arbitrary shape using $HP$-finite infinite elements |
| 96-20 | C. Schwab, M. Suri, C. Xenophontos | The $hp$ finite element method for problems in mechanics with boundary layers |
| 96-19 | C. Lage | The Application of Object Oriented Methods to Boundary Elements |
| 96-18 | R. Sperb | An alternative to Ewald sums. Part I: Identities for sums |
| 96-17 | M.D. Buhmann, Ch.A. Micchelli, A. Ron | Asymptotically Optimal Approximation and Numerical Solutions of Differential Equations |
| 96-16 | M.D. Buhmann, R. Fletcher | M.J.D. Powell's work in univariate and multivariate approximation theory and his contribution to optimization |
| 96-15 | W. Gautschi, J. Waldvogel | Contour Plots of Analytic Functions |
| 96-14 | R. Resch, F. Stenger, J. Waldvogel | Functional Equations Related to the Iteration of Functions |
| 96-13 | H. Forrer | Second Order Accurate Boundary Treatment for Cartesian Grid Methods |