# Séminaire Équations aux dérivées partielles – École Polytechnique

F. Treves

**Solution of Cauchy problems modulo flat functions**

*Séminaire Équations aux dérivées partielles (Polytechnique)* (1974-1975), exp. nᵒ 11, p. 1-19

<http://www.numdam.org/item?id=SEDP_1974-1975____A10_0>

S E M I N A I R E   G O U L A O U I C - L I O N S - S C H W A R T Z

1 9 7 4 - 1 9 7 5

## SOLUTION OF CAUCHY PROBLEMS MODULO FLAT FUNCTIONS

par F. TREVES

# § 1. INTRODUCTION

In the microlocal study of pseudodifferential operators (cf.[2]) one often encounters initial value problems, bearing on operator-valued functions of time, of the kind

(1) $\qquad \dfrac{\partial U}{\partial t} = i A(x,t,D_x)U \quad \underline{for} \quad 0 < t < T \ , \ U\big|_{t=0} = I$

(sometimes in view of forming solutions of

(2) $\qquad \dfrac{\partial E}{\partial t} - iA(x,t,D_x)E = I \quad , \quad$ the identity).

For the sake of simplicity let us assume, here, that

(3) $\qquad A(x,t,\xi) = \displaystyle\sum_{\nu=0}^{\infty} a_\nu(x,t,\xi) \ ,$

where $a_\nu(x,t,\xi)$, for each $\nu$, is a $C^\infty$ function of $(x,t,\xi)$, positive-homogeneous with respect to $\xi$ of degree $1 - \nu$, when x varies in an open set $\Omega$ of $\mathbb{R}^n$ and t in an interval $]-T,T[$ (though (1) is relative to the positive half-interval $[0,T[$). In some instances one is content with an approximate solution of (1) modulo operators which are regularizing (in the x-variables, and depend smoothly on t), sought in the form of a Fourier Integral Operator

(4) $\qquad U(t)u(x) = (2\pi)^{-n} \displaystyle\int e^{i\varphi(x,t,\xi)} k(x,t,\xi)\hat{u}(\xi)\,d\xi \ .$

Formally the __phase__ $\varphi$ should be determined by the __eikonal equation__ :

(5) $\qquad \varphi_t = a_0(x,t,\varphi_x) \ ,$

usually under the initial condition :

(6) $\qquad \varphi\big|_{t=0} = x \cdot \xi \ ,$

whereas the amplitude $k(x,t,\xi) = \sum\limits_{\nu=0}^{\infty} k_\nu(x,t,\xi)$ should satisfy the transport equations :

(7)
$$\frac{\partial k_\nu}{\partial t} - \sum_{j=0}^{\infty} \frac{\partial a_0}{\partial \xi_j}(x,t,\varphi_x)\frac{\partial k_\nu}{\partial x^j} + C(x,t;\varphi)k_\nu = F_\nu(x,t;\varphi;k_0,\ldots,k_{\nu-1}),$$

where $C(x,t;\varphi)$ and $F_\nu(x,t;\varphi;k_0,\ldots,k_{\nu-1})$ are the standard functionals ($F_\nu \equiv 0$ if $\nu = 0$). Furthermore the $k_\nu$ are usually submitted to the initial conditions :

(8)    __At time__ $t = 0$ , $k_0 \equiv 1$,    $k_\nu \equiv 0$   if $\nu > 0$.

The trouble with this method is that $a_0(x,t,\xi)$ need not be real, and thus (5)-(6) and (7)-(8) might not make sense (indeed, in general, the solution $\varphi$ should then be non real and what is then the meaning of $a_0(x,t,\varphi_x)$). But even if they make sense, e.g. when the $a_\nu(x,t,\xi)$ are (uniformly) analytic with respect to $\xi$, those Cauchy problems might then not be solvable.

Yet the method can be redeemed - at least when there exists a continuous function $\rho(x,t,\xi)$ with the properties we now describe. For convenience we shall assume $\rho$ positive-homogeneous of degree zero with respect to $\xi$. We shall solve (approximately) the above Cauchy problems on the unit sphere $|\xi| = 1$, and afterwards reestablish the appropriate homogeneity degrees. First of all we consider an extension of $a_0$ (and eventually of each $a_\nu$) to complex values $\xi + i\eta$ of $\xi$ of the form

(9)
$$\tilde{a}_0(x,t.\xi + i\eta) = \sum_{\alpha \in \mathbb{Z}_+^n} \frac{(i\eta)^\alpha}{\alpha!} g_\alpha(\eta) a_0^{(\alpha)}(x,t,\xi),$$

where $g_\alpha(\eta) = 1$ if $|\eta| \leq \varepsilon_\alpha$, $g_\alpha(\eta) = 0$ if $|\eta| > 2\varepsilon_\alpha$ ; also

(10)
$$\forall \lambda \in \mathbb{Z}_+^n, \quad |\partial_\eta^\lambda g_\alpha| \leq C_\lambda(1 + \varepsilon_\alpha^{-1})^{|\lambda|},$$

with $C_\lambda$ independent of $\varepsilon_\alpha > 0$. It is well known that we may choose the sequence $\varepsilon_\alpha \searrow +0$ such that the series at the right in (9) converges in $t^D(\tilde{a}_0$ is an almost-analytic extension of $a_0$ ; if $\tilde{\tilde{a}}_0$ is another extension,

analoguous to (9), $\widetilde{a}_o - \widetilde{\widetilde{a}}_o$ vanishes of infinite order at $\eta = 0$).

Suppose that, after substitution of the $\widetilde{a}_\nu$ for $a_\nu$ ($\nu = 0, 1, \ldots$) in (5) and in (7) (the $a_\nu$ for $\nu > 0$ enter in the definition of $C(x, t; \varphi)$ and of the $F_\nu$) we have found smooth solutions $\varphi, k_\nu$ of (5) and (7) - not exact ones, but only modulo functions which vanish of infinite order with respect to $\rho$ (more precisely, $\rho$-flat : see Def. 1) - and which also satisfy (6) and (8). Suppose furthermore that we prove that

$$(11) \qquad |Im\ \varphi_x| \leq C\rho\ , \quad \rho^d \leq C\ Im\ \varphi\ (|\xi| = 1)$$

for some constants $C \geq 0$, $d > 0$. Then the error arising from taking (4) as a solution of (1) will be of the form

$$R(t)u(x) = (2\pi)^{-n} \int e^{i\varphi(x, t, \xi)} r(x, t, \xi) \hat{u}(\xi)\, d\xi,$$

where, for any $N \in \mathbb{Z}_+$ and a suitable $C_N > 0$,

$$|r| \leq C_N |\xi| \rho^{Nd} \leq C_N' |\xi|^{1-N} (Im\ \varphi)^N,$$

and similar estimates hold for all the derivatives of r (we have exploited the homogeneities of the ingredients). A consequence of such inequalities is that R(t) is regularizing. It also follows from (11) and the preceeding remarks that different choices of the extensions $\widetilde{a}_o$, $\widetilde{\widetilde{a}}_1, \ldots$, lead to solutions (4) which only differ by regularizing operators.

In what precedes lies the motivation for the result presented here. Its generality far exceeds that required for solving (1) or (2). We hope that it will be an asset in future applications. Even when studying (1) and (2), one must make sure that the application of Th. 1 below leads to the existence of a function $\rho$ with the above properties, in particular (11). It is easy to see that this is not always so (otherwise all PDO's of principal type would have parametrices !). One still must investigate the relevant properties of the principal symbol $a_o(x, t, \xi)$. At the end we return briefly to this question.

## § 2. STATEMENT OF THE THEOREM

Definition 1 : Let X be a $C^\infty$ manifold, $\rho$ a continuous non negative function in X. We say that $f \in C^\infty(X)$ is $\rho$-flat if, given any integer N and any differential operator P (with $C^\infty$ coefficients) in X, $\rho^{-N}Pf$ is continuous in X.

If f is $\rho$-flat, we write $f \underset{\rho}{\sim} 0$ (if $f - g \underset{\rho}{\sim} 0$, we write $f \underset{\rho}{\sim} g$ , etc.).

In the sequel $\Omega$ denotes an open subset $\mathbb{R}^n$ (variable : $x = (x_1, \ldots, x_n)$), $u_0$ a $C^\infty$ mapping $\Omega \to \mathbb{R}^m$ ($m \geq 1$), $\mathcal{O}$ an open subset of $\mathbb{R}^{m(n+1)}$ which contains the image of $\Omega$ under the mapping :

$$(12) \qquad x \longmapsto (u_0(x), u_{0x}(x))$$

($f_x$ denotes the gradient of f ; subscripts mean differentiations). We are also given a $C^\infty$ mapping :

$$(13) \qquad F : \Omega \times \,]-T, T[\, \times \mathcal{O} \to \mathbb{R}^m \qquad (T > 0).$$

We study the (nonlinear) Cauchy problem :

$$(14) \qquad u_t = F(x, t, u, u_x), \quad x \in U , \quad |t| < \delta,$$

$$(15) \qquad u(x, 0) = u_0(x) , \quad x \in U ,$$

where $U \subset \Omega$ and $0 < \delta \leq T$ (we could as well limit the variation of t in (14) to $]0, \delta[$). With (14)-(15) we associate the ordinary diff. equ. (in which x plays the role of a parameter) :

$$(16) \quad v_t = F(x, t, v, q) , \quad q_t = F_x(x, t, v, q) + F_u(x, t, v, q)q ,$$

with initial conditions :

$$(17) \qquad v(x, 0) = u_0(x) , \quad q(x, 0) = u_{0x}(x) .$$

In (16), and throughout the forthcoming, we manipulate products of tensors with varying degrees of covariance and contravariance, according to the standard index contraction rule, without ever specifying what these degrees are. We note that, given any $U \subset\subset \Omega$, there is $\delta_U \in \, ]0,T]$ such that (16)-(17) has a unique solution $(v,q)$, $C^\infty$ in $\overline{U} \times [-\delta_U, \delta_U]$. For $(x,t)$ in the latter set we write :

(18)
$$\rho(x,t) = \left| \int_o^t |F_p(x,s,v(x,s),q(x,s))| \, ds \right|$$

$$(F_p(x,t,u,p) = \frac{\partial F}{\partial p}(x,t,u,p)).$$

**Theorem 1** : **Given any relatively compact open subset U of** $\Omega$ **there is a number** $\delta$, $0 < \delta \le \delta_U$, **such that the following is true** :

(I)     **There is** a $C^\infty$ **function** $u : U \times \, ]-\delta,\delta[ \, \rightarrow \, \mathbf{R}^m$ **verifying**

(19)
$$u_t - F(x,t,u,u_x) \underset{\rho}{\sim} 0 \text{ in } U \times \, ]-\delta,\delta[,$$

(20)
$$u(x,0) = u_o(x) \text{ in } U \, .$$

**Furthermore, there is a continuous function** $C(x,t) > 0$ **in** $U \times \, ]-\delta,\delta[$ **such that, in this set,**

(21)
$$\left| u - v \right| \le C\rho^2 \, , \quad \left| u_x - q \right| \le C\rho \, .$$

(II) **Any two** $C^\infty$ **mappings of** $U \times \, ]-\delta,\delta[$ **into** $\mathbf{R}^m$ **satisfying** (19)-(20) **are** $\rho$-**equivalent in** $U \times \, ]-\delta,\delta[$ .

By modifying $v$ and $q$ one can obtain more precise approximations of $u$ and of its derivatives (of any order) as shown in [3], sect. 4. Th.1 strengthens the main result of [3].

The proof of Th.1 is mainly that of the existence of the approximate solution $u$. Its uniqueness (mod. $\rho$-equivalence) will follow at once from one of the lemmas, and estimate (21) will also be a byproduct.

We begin by a standard <u>quasilinearization</u>, setting $w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix} = \begin{pmatrix} u \\ u_x \end{pmatrix}$

and transformaing (14)-(15) into

$$(22) \qquad w_t = A(x,t,w) + B(x,t,w)w_x, \qquad w(x,0) = \begin{pmatrix} u_o(x) \\ u_{ox}(x) \end{pmatrix} ,$$

where

$$A(x,t,w) = \begin{pmatrix} F(x,t,w_1,w_2) \\ F_x(x,t,w_1,w_2) + F_u(x,t,w_1)w_2 \end{pmatrix} , \quad B(x,t,w) = \begin{pmatrix} 0 & 0 \\ 0 & F_p(x,t,w_1,w_2) \end{pmatrix}$$

We call $\widetilde{w}$ the "vector" $\begin{pmatrix} v \\ q \end{pmatrix}$ in (16)-(17) ; the latter reads

$$(23) \qquad \widetilde{w}_t = A(x,t,w) \qquad , \qquad \widetilde{w}(x,0) = \begin{pmatrix} u_o(x) \\ u_{ox}(x) \end{pmatrix}$$

and we have

$$(24) \qquad \rho(x,t) = \left| \int_o^t |B(x,s,\widetilde{w}(x,s))| \, ds \right| .$$

## § 3. <u>PROOF OF THEOREM 1</u>

The first step consists in proving a weaker version of
Assertion (I), namely that Eq. (22) has an approximate solution, in a
suitable <u>fixed</u> interval $|t| < \delta$, modulo arbitrarily high, but finite,
powers of $|x|$ (i.e. functions vanishing of arbitrarily high order at
$x = 0$) and also modulo $\rho(0,t)$-flat functions.

To do this we study, in the ring of formal power series in x
whose coefficients are $C^\infty$ functions of t, $|t| < \delta$, the Cauchy problem :

$$(25) \quad \psi_t = A(x,t,\psi) + \left[ B(x,t,\psi) - B(0,t,\psi_o(t)) \right]\psi_x, \quad \psi(x,0) = \varphi(x) .$$

We write $\psi(x,t) = \sum_{\alpha \in \mathbb{Z}_+^n} \psi_\alpha(t)x^\alpha$ , $\varphi(x) = \sum_{\alpha \in \mathbb{Z}_+^n} \varphi_\alpha x^\alpha$ . As a matter of fact,

it is convenient to call $\psi_\nu(t)$ (resp. $\varphi_\nu$) the "vector" with components

$\psi_\alpha(t)$ (resp. $\varphi_\alpha$) for $|\alpha| = \nu$ ; $\nu = 0,1,\dots$ . In order to determine the

$\psi_\nu$ we differentiate $\nu$ times (25) with respect to x and put x = 0 in the result ; we obtain, thus :

$$(26)_0 \qquad \psi_0' = A(0,t,\psi_0) , \qquad\qquad \psi_0(0) = \varphi_0 \qquad ;$$

$$(26)_1 \qquad \psi_1' = G_1(t,\psi_0) + H_1(t,\psi_0,\psi_1)\psi_1, \quad \psi_1(0) = \varphi_1 ;$$

$$(26)_{\nu>1} \qquad \psi_\nu' = G_\nu(t,\psi_0,\ldots,\psi_{\nu-1}) + H_\nu(t,\psi_0,\psi_1)\psi_\nu , \; \psi_\nu(0) = \varphi_\nu .$$

The functionals $G_\nu$, $H_\nu$ are easy to draw from (25). They are polynomials with respect to $\psi_1,\ldots,\psi_{\nu-1}$ and $\psi_1$ resp. ($\deg_{\psi_1} H_\nu \leqq 1$) with coefficients of the form $A^{(\alpha,\beta,\gamma)}(0,t,\psi_0(t))$, $B^{(\alpha,\beta,\gamma)}(0,t,\psi_0(t))$. For us the crucial fact is the <u>linearity</u> of $(26)_{\nu>1}$ with respect to $\psi_\nu$ : It enables us to select $\delta > 0$, $\delta \leq \delta_U$, such that $(26)_0$ - $(26)_1$ have unique, smooth solutions $\varphi_0$, $\varphi_1$ for $|t| \leq \delta$, and then solve recursively $(26)_{\nu>1}$ without further decreasing $\delta$. Let us denote by $\mathcal{P}_\nu$ the (real) vector space of vectors such as $\varphi_\nu$, and set $\mathcal{P} = \mathcal{P}_0 \oplus \mathcal{P}_1 \oplus \ldots$ ($\mathcal{P}$ can be regarded as a linear space of formal power series in the $x_1,\ldots,x_n$ with coefficients in $\mathbf{R}^{m(n+1)}$). By solving $(26)_\nu$ we have defined a $C^\infty$ mapping $\varphi \mapsto \mathcal{G}(t,\varphi)$, depending smoothly on t, $|t| \leq \delta$, of an open subset of $\mathcal{P}$ , $\mathcal{U}$ , into $\mathcal{P}$. Actually this open set is of the form $\mathcal{U} = \mathcal{U}_0 \oplus \mathcal{U}_1 \oplus \mathcal{P}_2 \oplus \ldots \oplus \mathcal{P}_\nu \oplus \ldots$, where $\mathcal{U}_0$ is a neighborhood of $\begin{pmatrix} u_0(0) \\ u_{0x}(0) \end{pmatrix}$ and diam $\mathcal{U}_1$ is small enough. If $\psi(t) = \mathcal{G}(t,\varphi)$ we have

$$(27) \qquad \psi_\nu(t) = \mathcal{G}_\nu(t,\varphi_0,\ldots,\varphi_\nu) , \qquad \nu = 0,1,\ldots .$$

Consider then the Fréchet derivative $\dfrac{D\psi}{D\varphi}$ (at some point of $\mathcal{U}$). In the direct sum decomposition $\mathcal{P} = \underset{\nu}{\oplus} \mathcal{P}_\nu$ it is represented by a triangular matrix, according to (27). It will be an automorphism of $\mathcal{P}$ if, for every $\nu$, $\dfrac{D\psi_\nu}{D\varphi_\nu}$ is an automorphism of $\mathcal{P}_\nu$ . Since $\mathcal{G}_\nu|_{t=0}$ = Identity of $\mathcal{P}_\nu$ . the last assertion is standard. Thus, we have proved

(28) $\quad \dfrac{D\mathscr{C}}{D\varphi}$ is an automorphism of $\mathscr{P}$, depending on $C^\infty$ fashion on $t, |t| \leq \delta$.

We introduce an (unknown) formal power series in x, $\varphi(t)$, depending smoothly on t, and set

(29) $$w(t) = \mathscr{C}(t, \varphi(t)).$$

We put this into (22') which will serve to determine $\varphi$ ; for $|t|$ small $\varphi_0(t)$ will stay in $\mathscr{U}_0$, $\varphi_1(t)$ in $\mathscr{U}_1$. Computation yields :

(30) $\quad \dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)\varphi_t = A(x,t,\mathscr{C}(t,\varphi)) - \mathscr{C}_t(t,\varphi) + B(x,t,\mathscr{C}(t,\varphi))\dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)\varphi_x)$

$$\varphi\big|_{t=0} = \begin{pmatrix} u_o \\ u_{ox} \end{pmatrix} ,$$

recalling that $\mathscr{C}(0,\varphi) = \varphi$. But, by (25) ,

(31) $\mathscr{C}_t(t,\varphi) = A(x,t,\mathscr{C}(t,\varphi)) + [B(x,t,\mathscr{C}(t,\varphi)) - B(0,t,\mathscr{C}_0(t,\varphi_0))]\dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)\varphi_x$ ,

since $\partial_x[\mathscr{C}(t,\varphi)] = \dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)\varphi_x$ . Putting (31) into (30) yields :

(32) $\quad \varphi_t = B^\sharp(t,\varphi)\varphi_x$ , $\qquad \varphi\big|_{t=0} = \begin{pmatrix} u_o \\ u_{ox} \end{pmatrix}$

where

(33) $\quad B^\sharp(t,\varphi) = \dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)^{-1}B(0,t,\mathscr{C}_0(t,\varphi_0))\dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)$ .

At this point we recall that $\mathscr{C}_0(t,\varphi_0) = \psi_0$, the solution of $(26)_0$.
But here $\psi_0(0) = \begin{pmatrix} u_o(o) \\ u_{ox}(o) \end{pmatrix}$ . Comparing with (23) shows that

$\mathscr{C}_0(t,\varphi_0) = \widetilde{w}(0,t)$, hence :

(34) $\qquad B^\sharp(t,\varphi) = \dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)^{-1}B(0,t,\widetilde{w}(0,t))\dfrac{D\mathscr{C}}{D\varphi}(t,\varphi)$ .

We apply an elementary lemma about ODES (in a Banach space with norm $\| \ \|$) .

Let $e_0 \in E$, $V$ be an open neighborhood of $e_0$, $g(t,e)$ a continuous function in $[-\delta, \delta] \times V$, valued in $E$, Lipschitz continuous with respect to $e$ :

$$(35) \qquad \| g(t,e_1) - g(t,e_2) \| \leq K \| e_1 - e_2 \|, \quad e_1, e_2 \in V, \quad |t| \leq \delta.$$

We introduce the function :

$$(36) \qquad \rho_0(t) = \left| \int_0^t \| g(s,e_0) \| ds \right|.$$

**Lemma 1** : Under the preceeding hypotheses there is a number $\varepsilon > 0$ and a $C^1$ function $\Phi$ of $t$, $|t| \leq \delta$, valued in $V \subset E$, such that

$$(37) \qquad \Phi_t = g(t,\Phi) \quad \text{in the set} \quad |t| < \delta, \ \rho_0(t) < \varepsilon,$$

$$\Phi(0) = e_0.$$

**Moreover,** if $|t| < \delta$ and $\rho_0(t) < \varepsilon$,

$$(38) \qquad \left\| \Phi(t) - e_0 - \int_0^t g(s,e_0) ds \right\| \leq (e^{K|t|} - 1) \rho_0(t).$$

**Proof** : Apply Picard's method, taking $\Phi = \Phi^1 + \sum_{j=1}^{\infty} (\Phi^{j+1} - \Phi^j)$ where

$$\Phi^j(t) = e_0 + \int_0^t g(s, \Phi^{j-1}(s)) ds \quad (\Phi^0 \equiv e_0). \text{ We have}$$

$$(39) \qquad \| \Phi^{j+1}(t) - \Phi^j(t) \| \leq (K|t|)^j \rho_0(t) / j! , \quad j = 0, 1, \ldots .$$

In particular, $\| \Phi^j(t) - e_0 \| \leq \varepsilon \, e^{K\delta}$ if $|t| \leq \delta$, $\rho_0(t) \leq \varepsilon$, and thus $\Phi^j(t) \in V$ if $\varepsilon > 0$ is small enough. The conclusion of Lemma 1 follows at once.

Of course, if $g(t,e)$ is smooth, so is the solution $\Phi$. We apply Lemma 1 to problem (32) - taking $V$ to be the space of polynomials with respect to $x$ of degree $\leq N$ (with coefficients in $\mathbb{R}^{m(n+1)}$). Of course in doing so we introduce an error, which is $O(|x|^{N+1})$. We are faced with a problem

$$(40) \qquad \partial_t \Phi_N = F_N(t, \Phi_N) \;, \quad \Phi_N \big|_{t=0} = \sum_{|\alpha| \le N} \frac{x^\alpha}{\alpha!} \partial_x^\alpha \varphi \big|_{t=0} \;.$$

From (34) (and (28)) it follows at once that

$$(41) \qquad \| F_N(t, \Phi_N(0)) \| \le C_N |B(0, t, \widetilde{w}(0, t))| \;.$$

Thus we may take $\rho_o(t)$ in lemma 1 equal to $C_N \rho(0, t)$ (see (24)).

We set

$$(42) \qquad w_N(t) = \mathcal{C}(t, \Phi_N(t)) \;;$$

by (27) we know that $w_N$ is a polynomial with respect to x of degree $\le N$ ; we regard $w_N(t)$ as a (smooth) function of $(x, t)$. It verifies, as we see by retracing our steps through (30)

$$(43) \qquad \partial_t w_N = A(x, t, w_N) + B(x, t, w_N) \partial_x w_N + O(|x|^{N+1})$$

$$\text{if } |t| \le \delta \;, \; \rho(0, t) \le \varepsilon_N \;,$$

$$(44), \qquad w_N(x, 0) - \begin{pmatrix} u_o(x) \\ u_{ox}(x) \end{pmatrix} = O(|x|^{N+1}) \;.$$

This implies at once our claim, at the beginning of § 3. Note that if the functions A, B and the initial datum $u_o$ are $C^\infty$ functions of a parameter z, and if the numbers $\delta_U$, $\delta$ can be chosen independently of z (which is true if z varies in a compact manifold), the approximate solutions $w_N$ can be   selected so as to depend smoothly on z (and $\delta$ does not have to be decreased through the argument).

## § 4.   END OF THE PROOF OF THEOREM 1

At this stage we connect the argument to that in pp. 352-367 of [3]. We apply the result in sect. 3 to the problem

$$(45) \qquad w_t = A(x+y, t, w) + B(x+y, t, w) w_y \;, \quad w \big|_{t=\cdot} = \begin{pmatrix} u_o(x+y) \\ u_{ox}(x+y) \end{pmatrix} \;.$$

Here x plays the role of a    parameter, varying in $U \Subset \Omega$. The number $\delta$
of sect. 3 can be chosen independently of x. The analogue of $\tilde{w}(x,t)$
here is $\tilde{w}(x+y,t)$ and that of $\rho(x,t)$ is $\rho(x+y,t)$. This means that the
role of $\rho(0,t)$ will now be played by $\rho(x,t)$. We denote by $W_N(x,y,t)$ the
corresponding approximate solutions of (45), mod. $O(|y|^{N+1})$, in a set

$$(46) \qquad \Sigma_N = \{(x,t) \; ; \; x \in U, \; |t| < \delta , \; \rho(x,t) < \varepsilon_N\} \qquad (N = 0,1,\ldots).$$

Let us set $\theta = \theta(x,t) = w_N(x,0,t) - \tilde{w}(x,t)$. We have :

$$(47) \qquad \theta_t = A(x,t,\tilde{w}+\theta) - A(x,t,\tilde{w}) + B(x,t,\tilde{w}+\theta)\partial_y w_N(x,0,t), \quad \theta\big|_{t=0} = 0 .$$

We apply the last part of lemma 1 to (47). We get :

$(48) \qquad \underline{\text{if } \varepsilon_N > 0 \text{ is small enough and } C_N > 0 \text{ large enough,}}$

$$|w_N(x,0,t) - \tilde{w}(x,t)| \le C_N \rho(x,t) \text{ in } \Sigma_N .$$

We consider the approximate eq. (45) (where $w_N$ replaces $w$)
and differentiate with respect first to x, then to y, and subtract :

$$(49) \qquad \partial_t(w_{N,x} - w_{N,y}) =$$

$$\{A_w(x+y,t,w_N) + B_w(x+y,t,w_N)\partial_y w_N\} \times (w_{N,x} - w_{N,y}) + B(x+y,t,w_N)\partial_y(w_{N,x} - w_{N,y}) + O(|y|^N)$$

$$(50) \qquad w_{N,x} - w_{N,y}\big|_{t=0} = O(|y|^N) .$$

We set, for any $\alpha,\beta \in \mathbb{Z}_+^n$ , $\ell \in \mathbb{Z}_+$ , $|\beta| < N$ ,

$$\theta_{\alpha,\beta,\ell} = \partial_x^\alpha \partial_y^\beta \partial_t^\ell (w_{N,x} - w_{N,y})\big|_{y=0}$$

We apply $\partial_x^\alpha \partial_y^\beta \partial_t^\ell$ to both sides of (49) and make y = 0 in the result ,
obtaining :

$$(51) \qquad \partial_t \theta_{\alpha,\beta,\ell} = \Sigma^* C_{\alpha,\beta,\ell}^{\alpha',\beta',\ell'}(x,t)\theta_{\alpha',\beta',\ell'} +$$

$$+ B(x,t,w_N(x,0,t))\{\partial_y[\partial_x^\alpha \partial_y^\beta \partial_t^\ell (w_{N,x} - w_{N,y})]\}_{y=0} .$$

The summation in $\Sigma^*$ extends to triples $\alpha',\beta',\ell'$ such that $|\alpha'+\beta'| + \ell' \leq |\alpha+\beta| + \ell$, $\alpha'_j \leq \alpha_j$ $(j = 1,\ldots,n)$, $\ell' \leq \ell$ and $|\beta'| \leq |\beta| + 1$.

From (50) and the fact that $|\beta| < N$ we derive

$$(52) \qquad \theta_{\alpha,\beta,0}\big|_{t=0} = 0$$

Putting this into (51) and induction on $(\beta,\ell)$, $|\beta| + \ell < N$, shows that

$$(53) \qquad \theta_{\alpha,\beta,\ell}\big|_{t=0} = 0 \qquad \text{if } |\beta| + \ell < N.$$

We apply Lemma 1 to problem (51)-(53), regarding the $\theta_{\alpha,\beta,\ell}$ for $|\alpha + \beta| + \ell \leq k$, as unknowns, while $\theta_{\alpha,\beta,\ell}$, $|\alpha+\beta|+\ell = k+1$, play the role of data (this for $k = 0,1,\ldots,N-1$). We get :

$$(54) \qquad \sum_{|\alpha+\beta|+\ell = k} |\theta_{\alpha,\beta,\ell}(x,t)| \leq$$

$$C'_N(x,t) \left| \int_o^t |B(x,s,w_N(x,0,s))| \, ds \right| \sup_{s\in[0,t]} \sum_{|\alpha+\beta|+\ell = k+1} |\theta_{\alpha,\beta,\ell}(x,s)| \ .$$

Here $C'_N$ is a continuous positive function in $\Sigma_N$.

At this stage we take advantage of (48) ; it implies (in $\Sigma_N$)

$$(55) \qquad \left| \int_o^t |B(x,s,w_N(x,0,s)| \, ds \right| \leq CC_N \rho(x,t),$$

hence, by (54) ,

$$(56) \qquad \sum_{|\alpha+\beta|+\ell = k} |\theta_{\alpha,\beta,\ell}(x,t)| \leq C''_N \left\{ \sup_{s\in[0,t]} \sum_{|\alpha+\beta|+\ell = k+1} |\theta_{\alpha,\beta,\ell}(x,s)| \right\} \rho(x,t).$$

By descending induction on $k = N-1,\ldots,0$, we obtain, in $\Sigma_N$,

$$(57) \qquad \sum_{|\alpha+\beta|+\ell = k} |\theta_{\alpha,\beta,\ell}(x,t)| \leq C_N^{(iii)}(x,t)\rho(x,t)^{N-k} \ .$$

Actually we shall need the following consequence of (57) :

(58)
$$\sum_{|\alpha|+\ell=k} |\theta_{\alpha,0,\ell}| \le C_N^{(iii)} \rho^{N-k} \quad \underline{in} \quad \Sigma_N \; .$$

From now we set $\omega_N(x,t) = w_N(x,0,t)$. We have :

(59)
$$\partial_t \omega_N - A(x,t,\omega_N) - B(x,t,\omega_N)\partial_x \omega_N = R_N(x,t) \; ,$$

(60)
$$\omega_N(x,0) = \begin{pmatrix} u_o(x) \\ u_{ox}(x) \end{pmatrix}$$

where

(61)
$$R_N(x,t) = B(x,t,\omega_N)\left[w_{N,y}(x,0,t) - w_{N,x}(x,0,t)\right]$$

satisfies, for all $k = 0,1,\ldots,N$ (according to (58))

(62)
$$\sum_{|\alpha|+\ell=k} |\partial_x^\alpha \partial_t^\ell R| \le C_N^{(iv)} \rho^{N-k} \quad \underline{in} \; \Sigma_N \; .$$

$(C_N^{(iv)}$ is a continuous positive function in $\Sigma_N)$. We shall apply the following

<u>Lemma 2</u> : <u>Let</u> $\chi^j$ $(j = 1,2)$ <u>be any two</u> $C^\infty$ <u>functions in</u> $\Sigma_N$, <u>satisfying</u> <u>there</u>

(63)
$$\partial_t \chi^j - A(x,t,\chi^j) - B(x,t,\chi^j)\partial_x \chi^j = R^j(x,t) \; ,$$

(64)
$$\chi^j(x,0) = \begin{pmatrix} u_o(x) \\ u_{ox}(x) \end{pmatrix} \; ,$$

<u>where the</u> $R^j$ $(j = 1,2)$ <u>verify</u>

(65)
$$R^j/\rho \in C^o(\Sigma_N) \; ,$$

(66) <u>for every</u> $k = 0,1,\ldots,N$, $\rho^{-N+k} \sum_{|\alpha|+\ell=k} |\partial_x^\alpha \partial_t^\ell (R^1 - R^2)| \in C^o(\Sigma_N) \; .$

**Then, for every** $k = 0, 1, \ldots, N$ ,

(67)
$$\rho^{-N+k} \sum_{|\alpha|+\ell-k} |\partial_x^\alpha \partial_t^\ell (\chi^1 - \chi^2)| \in C^0(\Sigma_N) \quad .$$

**Proof :** Set $\chi = \chi^1 - \chi^2$, $R = R^1 - R^2$, whence, by (63) - (64),

(68)
$$\chi_t = A(x,t,\chi+\chi^2) - A(x,t,\chi^2) + B(x,t,\chi+\chi^2) \partial_x \chi +$$

$$+ [B(x,t,\chi+\chi^2) - B(x,t,\chi^2)] \partial_x \chi^2 + R ,$$

and $\chi(x,0) = 0$. Let us rewrite (68) in the form

(69)
$$\chi_t = \mathcal{U}(x,t)\chi + B(x,t,\chi^2)\chi_x + R$$

and apply $\partial_x^\alpha \partial_t^\ell$ $(|\alpha| + \ell = k < N)$ to both sides of (69) .

We obtain :

(70)
$$\partial_t(\partial_x^\alpha \partial_t^\ell \chi) = \sum_{\substack{|\alpha'|+\ell' \leq k \\ \ell' \leq \ell}} C_{\alpha,\ell}^{\alpha',\ell'}(x,t) \partial_x^{\alpha'} \partial_t^{\ell'} \chi +$$

$$+ B(x,t,\chi^2)\partial_x(\partial_x^\alpha \partial_t^\ell \chi) + \partial_x^\alpha \partial_t^\ell R \quad .$$

Induction on $\ell = 0, 1, \ldots, N-1$ , shows that

(71)
$$\partial_x^\alpha \partial_t^\ell \chi(x.0) = 0 \quad .$$

We regard (70)-(71) as an initial value problem for a system of **linear** ODES, with unknowns the $\partial_x^\alpha \partial_t^\ell \chi$, $|\alpha|+\ell \leq k$, and data the same but where $|\alpha|+\ell = k+1$, and we let $k$ range over $0,1,\ldots,N-1$. We apply Lemma 1 ; because of the linearity of our system with respect to the unknowns we are not forced to decrease the number $\varepsilon_N$ in (46). We get, for a suitable positive function $\tilde{C}_N \in C^0(\Sigma_N)$,

(72)
$$\sum_{|\alpha|+\ell=k} |\partial_x^\alpha \partial_t^\ell \chi(x,t)| \leq \tilde{C}_N(x,t)\rho(x,t)^{N-k} +$$

$$+ \left| \int_0^t |B(x,s,\chi^2(x,s))| \, ds \right| \sup_{s \in [0,t]} \left\{ \sum_{|\alpha|+\ell=k+1} |\partial_x^\alpha \partial_t^\ell \chi(x,t)| \right\} \quad .$$

We have taken (66) into account. We may reach the desired conclusion if we reason by descending induction on k and if we show that

$$(73) \qquad |\int_{o}^{t} |B(x,s,\chi^2(x,s))| \, ds| \leq const. \ \rho(x,t),$$

which in turns follows from the fact that

$$(74) \qquad |\chi^2 - \tilde{w}| \ / \rho \ \in \ C^0(\Sigma_N) \quad .$$

In order to prove (74) it suffices to observe that $\eta = \chi^2 - \tilde{w}$ satisfies

$$(75) \qquad \eta_t \ = \ \mathcal{Q}\ell^\#(x,t)\eta + B(x,t,\tilde{w})\chi_x^2 + R^2(x,t), \qquad \eta(x,0) \ = \ 0 \ .$$

We apply once more Lemma 1 regarding (75) as an initial value problem for a linear ODE with unknown $\eta$) and we get easily (74).

We can conclude the proof of Th. 1. First of all notice that by applying Lemma 2 with $N \nearrow +\infty$ we get at once the uniqueness part, (II). On the other hand, returning to (59)-(60) and applying Lemma 2 with $\chi^j = \omega_{N+j-2}$ (j = 1,2 ; we assume N > 1 and $\varepsilon_N \searrow +0$ as $N \nearrow +\infty$ ). We see that, for a suitable constant $M_N > 0$ and every $k = 0,1,\ldots,N-1$,

$$(76) \qquad \rho^{N-1-k} \overline{\sum_{|\alpha|+\ell=k}} |\partial_x^\alpha \partial_t^\ell (\omega_N - \omega_{N-1})| \ \leq \ M_N \quad \underline{in} \quad \Sigma_N \quad ,$$

By using a partition of unity à la Whitney we may find a sequence $\{\zeta_N\}$ of $C^\infty$ functions in $U \times ]-\delta,\delta[$, having the following properties :

$$(77) \qquad \zeta_N(x,t) \ = \ 1 \quad \underline{if} \quad \rho(x,t) < \varepsilon_N'/2 \ , \quad \zeta_N(x,t) \ = \ 0 \ \underline{if} \ \rho(\cdot,t) \geq \varepsilon_N' \ ;$$

$$(78) \qquad \forall \alpha \in \mathbb{Z}_+^n \ , \quad \forall \ell \in \mathbb{Z} \ , \quad \exists \ C_{\alpha,\ell} > 0 \quad \underline{such \ that}, \ \underline{for \ all \ N \ and \ all}$$

$$x \in U, \ t \in ]-\delta,\delta[ \ ,$$

$$|\partial_x^\alpha \partial_t^\ell \zeta_N(x,t)| \ \leq \ C_{\alpha,\ell} (1 + 1/[\varepsilon_N' \rho(x,t)])^{|\alpha|+\ell} \quad .$$

We shall take $\varepsilon_N' \leq \varepsilon_{N+1}$ in such a way that the series at the right, in

$$(79) \qquad w = \omega_1 + \sum_{N=1}^{+\infty} \zeta_N (\omega_{N+1} - \omega_N) ,$$

converges in $C^\infty(U \times ]-\delta, \delta[)$, and that w satisfies (45) if we interpret the first eq. (45) as a congruence mod $\rho$-flat functions. By taking advantage of (76), it is easy to see that such a choice is possible.

Let us remark that if $w = \begin{pmatrix} w_1 \\ w_2 \end{pmatrix}$ is the approximate solution (79) and if we set $u = w_1$, we have $u_t - F(x,t,u,w_2) \underset{\rho}{\sim} 0$, whence, by differentiating this with respect to x, and taking (22) into account :

$$(80) \quad \partial_t(u_x - w_2) - F_u(x,t,u,w_2)(u_x - w_2) \underset{\rho}{\sim} 0, \quad (u_x - w_2)\big|_{t=0} = 0 ,$$

whence immediately $u_x \underset{\rho}{\sim} w_2$ and, consequently, (19). By subtracting (24) from (22) we easily see that $(w - \tilde{w})/\rho$ is continuous in $U \times ]-\delta, \delta[$ which implies that $\rho^{-1}\{|u-v| + |u_x - q|\}$ is also continuous there. But we have also (cf. (16) and (19)) :

$$(81) \qquad (u-v)_t \underset{\rho}{\sim} F(x,t,v+(u-v),q) - F(x,t,v,q) + F_p(x,t,v+(u-v),q)(u_x - q)$$

$$+ O(|u_x - q|^2) ,$$

with $(u-v)\big|_{t=0} = 0$ . We apply once more the estimate (38). We obtain that, modulo $\rho$-flat functions ,

$$|(u-v)(t)| \leq C\{ |\int_0^t |F_p(x,s,v(x,s)||(u_x - q)(x,s)|ds|$$

$$+ |(u_x - q)(x,t)|^2\} \leq C_1 \rho(x,t)^2 ,$$

in view of what we have just said (C and $C_1$ are positive continuous functions of $(x,t)$).

Q.E.D.

## § 5. RETURN TO THE EIKONAL AND TRANSPORT EQUATIONS

We apply Th. 1 to the problems (5)-(6) and (7)-(8). We may assume that $a_o = \sqrt{-1}\ b_o$ with $b_o$ real. Indeed, by means of a canonical transformation (which only affects $U(t)$ by changing it into $V(t)^{-1} U(t) V(t)$, where $V(t)$ is an elliptic F.I.O. with real phase) we may reduce the situation to this case. In order that the errors resulting from application of Th. 1 be negligeable (cf.(11)) we need some hypotheses on $b_o$. We suppose that

$$(82) \qquad b_o(x,t,\xi) \geq 0 \qquad \underline{\text{for all }} (x,\xi) \ \underline{\text{and}}\ t \geq 0$$

(actually the argument is local). It is easy to see that (82) implies, in the same region, with a suitable $M > 0$,

$$(83) \qquad |\text{grad}_{x,\xi}\ b_o| \leq M \sqrt{b_o}.$$

We apply Th. 1 to (5)-(6) where $a_o = ib_o$ has been replaced by $\tilde{a}_o = i\tilde{b}_o$ given by (9). Of course we view (5)-(6) as bearing on functions valued in $\mathbb{R}^2$, but it is convenient to maintain the identification $\mathbb{R}^2 \cong \mathbb{C}$. The particular function $F(x,t,u,p)$ in this case is independent of $u$, and the defining equations (16)-(17) of $(v,q)$ simplify. First of all,

$$(84) \qquad q_t = f_x(x,t,q)\ , \qquad q(x,0) = (\xi,0) \qquad (\in \mathbb{R}^{2n}).$$

Estimate (38) applied to (84) yields :

$$(85) \qquad |q(x,t,\xi) - (\xi,0)| \leq \text{const.}\left|\int_o^t |F_x(x,s,\xi)|\,ds\right|,$$

whence by (83) ,

$$\rho(x,t,\xi) = \left|\int_o^t |F_p(x,s,q(x,s,\xi))\,ds\right| \leq \text{const.}\left|\int_o^t |\nabla_{x,\xi} F(x,s,\xi)|\,ds\right|$$

$$\leq \text{const.}\left|\int_o^t |\nabla_{x,\xi} b_o(x,s,\xi)|\,ds\right| \leq \text{const.}\left|\int_o^t \sqrt{b_o(x,s,\xi)}\,ds\right|.$$

Thus, if we set

(86)
$$B_0(x,t,\xi) = \int_0^t b_0(x,s,\xi)\,ds,$$

we get, for $t \geq 0$,

(87)
$$\rho(x,t,\xi) \leq \text{const.} \left[t\,B_0(x,t,\xi)\right]^{1/2} .$$

Next we look at the equations defining $v$ :

(88)
$$v_t = i\,\tilde{b}_0(x,t,q) , \qquad v(x,0) = x.\xi \qquad (\in \mathbb{R} \subset \mathbb{C}).$$

Again by (38), this yields

(89)
$$\left| v - x.\xi - iB_0 \right| \leq \text{const.} \; tB_0 .$$

By the first inequality (21) and since $\rho^2 \leq \text{const.} \; tB_0$ ,

(90)
$$\left| \varphi - x.\xi - iB_0 \right| \leq \text{const.} \; tB_0 .$$

we apply now the second inequality (21), in conjunction with (85) :

(91)
$$\left| \varphi_x - \xi \right| \leq \text{const.} \; (tB_0)^{1/2} .$$

This, i.e. (87), (90), (91), gives us all we want : the assertions in Sect. 1 concerning the eikonal equations are all valid if we substitute $tB_0$ for $\rho$ (and take the time interval $[0,T[$ short enough), in particular (11), with $d = 2$.

Finally we look at (7)-(8), where suitable extensions $\tilde{a}_\nu$ have been substituted for $a_\nu$. The relevant observation, here, is that when $F$ is _affine_ with respect to $(u,p)$, i.e.

(92)
$$F(x,t,u,p) = A(x,t)p + B(x,t)u + C(x,t),$$

the function denoted by $\rho$ in Th. 1 is nothing else but $\left| \int_0^t |A(x,s)|\,ds \right|$. In our case it will be

(93)
$$\rho_1(x,t,\xi) = \left| \int_0^t |\nabla_\xi\, b_0(x,s,\varphi_x(x,s,\xi))|\,ds \right| .$$

If we apply (91) we find at once that $\rho_1$ is also $\leq$ const. $(tB_0)$ (cf. proof of (87)). As a conclusion we see that the found phase and amplitude are approximate solutions modulo $\rho$-flat functions (with now $\rho = tB_0$) of (5) and (7) resp., and that (11) holds. Furthermore a careful inspection of the estimate for the length of the time interval $[0,\delta]$ would show that it depends only on (uniform) properties of $b_0(x,t,\xi)$, thus opening the door to globalization. On this we refer to works of Melin and Sjöstrand, in particular [1]. Let us also mention that the argument outlined in sect. 1 and 5 of the present article yields a construction of parametrices for the class of differential operators considered in [2], and also for wider classes of pseudodifferential operators.

---

## REFERENCES

[1]  Melin, A. and Sjöstrand, J. - Fourier integral operators with complex-valued phase functions, to appear.

[2]  Treves, F. - Hypoelliptic partial differential equations of principal type. Sufficient conditions and necessary conditions. Comm. Pure Applied Math., vol. XXIV (1971), 631-670.

[3]  Treves, F. - Approximate solutions to Cauchy problems, J. of Diff. Eq., Vol. 11 (1972), 349-363.

[4]  Treves, F. - Local solvability in $L^2$ of first order linear PDES, Amer. J. of Math., Vol XCII (1970), 369-380.

---