

PROPULSION

ERNEST YEUNG [ERNESTYALUMNI@GMAIL.COM](mailto:ERNESTYALUMNI@GMAIL.COM)

CONTENTS

Part 1. Notes and Solutions for *Rocket Propulsion Elements* by George Sutton and Oscar Biblarz

- 1. Definitions and Fundamentals
- 2. Nozzle Theory and Thermodynamic Relations
- 3. Flight Performance

Part 2. Notes and Solutions for *Space Propulsion Analysis and Design* by Humble, Henry, Larson

- 4. Thermodynamics of Fluid Flow; cf. Ch. 3. of Humble, Henry, and Larson (1995) [4]

Part 3. 1-dim. propulsion (revisited)

- 5. Fluid Flow (review)
- 6. Ideal Rocket Equation
- 7. Thermochemistry, Thermodynamics, One-dimensional Fluid Flow
- 8. Multiple Staging
- 9. Lagrangian
- 10. Mission Design

Part 4. AE121

- 11. Isentropic Flow Eqns. with Area Change
- 12. PSs
- 13. Equilibrium flow vs. frozen flow for Nozzle Flow and the Example of the Space Shuttle Main Engine (SSME)
- 14. Liquid-Vapor Equilibrium

Part 5. Basic Feeling

- 15. Box with a hole rocket; bottled (box) rocket

Part 6. Combustion

- 16. mass fraction vs. mole fraction vs. molecular mass i.e. “molecular weight”
- 17. Enthalpy
- 18. Thermochemistry of combustion
- 19. Droplet Evaporation
- 20. Droplet Evaporation and Burning
- 21. Shvab-Zeldovich forms
- 22. Droplet model; Burning Droplets
- 23. Combustion Chamber Flow Model

Part 7. Numerical Computation; Scientific Computation

*Date:* 13 novembre 2015.

*Key words and phrases.* Propulsion, Rocket Propulsion, Thermodynamics, Fluid Flow, Fluid Mechanics.

24.	Interpolation and Extrapolation	40
25.	Integration of ODEs: Runge-Kutta Methods	41
26.	Systems of Ordinary Differential Equations ODEs	51
27.	Computational Fluid Dynamics (CFD) and Navier-Stokes equation	51
28.	Hamiltonian	56
29.	Heat Equation	56
30.	Lattice Boltzmann method	56
31.	Reactive Flows	57
32.	Sparse Matrices	57
8	Part 8. Optimization	57
8	33. Minimization or Maximization of Functions	57
8	34. Conjugate Gradient Methods	60
12		
12	Part 9. Cantera	63
12	35. Professor David Goodwin, 1957-2012	63
13	36. Speed of Sound and sound speed.py	64
13		
14	Part 10. Compressible Flow	65
15	37. Nozzle Flow of a Reacting Gas	65
23	38. Dynamic pressure	67
25		
25	Part 11. Electromagnetism	67
25		
25	Part 12. Physical Kinetics	69
27		
27	Part 13. Guidance, Navigation, and Control (GNC)	69
28	39. Software for modeling and execution	69
28		
28	Part 14. Control, Control Theory	69
30	40. Signals and Systems	69
32	41. Introduction to Control Systems	69
32	42. Mathematical Modeling of Control Systems, Block Diagrams	69
34	43. Mathematical Modeling of Mechanical Systems	73
37	44. Discrete Control, s-domain to z-domain	74
	45. Spacecraft Attitude Control	74
40	46. More resources	75

Part 15. Orbital Mechanics

- 47. *N*-body problem
- 48. Orbit Determination from Observations; Coordinate Systems

Part 16. GPS, Geodesy

- 49. Geodesy

Part 17. Real-time systems

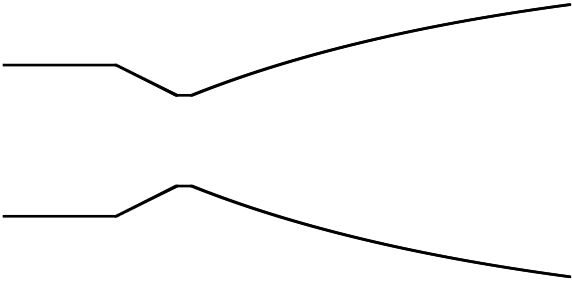
- 50. Partial Ordering of Events in a Distributed System
- 51. Logical clocks
- 52. Vector Clocks
- 53. Clock Synchronization, Lamport Timestamps, Vector Clocks, References

Part 18. Flight Software

References

ABSTRACT. Everything about Propulsion, with a focus on rocket propulsion.

I also look at (rocket) propulsion for engineers from a (theoretical and mathematical) physicists’ point of view. I would like to seek more cross-pollination between physicists and mathematicians and engineers in thermodynamics and fluid mechanics.



Part 1. Notes and Solutions for *Rocket Propulsion Elements* by George Sutton and Oscar Biblarz

1. DEFINITIONS AND FUNDAMENTALS

Ch. 2 Definitions and Fundamentals of Biblarz and Sutton (2001) [3]

Biblarz and Sutton (2001) [3] says that

Propulsion is achieved by applying a force to a vehicle;

This propulsive force is obtained by ejecting propellant at high velocity.

Let surface  $S$  of the rocket (with remaining fuel), choosing the normal unit vector of  $S$ ,  $\mathbf{n} \equiv \mathbf{n}_S$  to face outward.

Let the force, per unit area, field on a vehicle (also another word for rocket, plus remaining fuel)  $\mathbf{f} = \mathbf{f}(t, \mathbf{x})$ ,  $\forall \mathbf{x} \in S$ .

From pp. 144 of Frankel (2011) [6], vector integrals *make no sense* on general manifolds (how could we add 2 vectors located at different points?!). The moral is that we naturally integrate differential forms, not vectors (i.e. vector fields) (cf. pp. xxxviii of Frankel (2011) [6]).

75

75

77

$$\mathfrak{X}(S) \xrightarrow{\flat} \Omega^1(S) \xrightarrow{\quad \quad \quad *} \Omega^{d-1}(S)$$

77

77

$$\mathbf{f} = f^i \frac{\partial}{\partial x^i} \mapsto g_{ij} f^i dx^j = f_j dx^j \mapsto \frac{\sqrt{g}}{(d-1)!} f_j g^{jk} \epsilon_{k k_2 \dots k_d} dx^{k_2} \wedge \dots \wedge dx^{k_d} = f_j g^{jk} dS_k = f^i dS_i$$

77 for  $dS_i := \frac{\sqrt{g}}{(d-1)!} \epsilon_{ii_2 \dots i_d} dx^{i_2} \wedge \dots \wedge dx^{i_d} \in \Omega^{d-1}(S)$ ,

77 and if  $g_{ij} g^{jk} = \delta_i^k$

78 Note that  $f$  and  $S$  are time-dependent in a fixed, inertial frame:

79

$$\mathbf{f} = \mathbf{f}(t, \mathbf{x}), \quad \forall \mathbf{x} \in S$$

79

$$S = S(t)$$

Interpret

79

80

$$(1) \qquad \qquad \qquad \int_S *(\mathbf{f}(t, \mathbf{x}))^\flat \equiv F(t) \in \mathbf{R}$$

as the net amount of force exerted on a rocket (with remaining fuel), by expelled propellants (at that instance, at the surface  $S$ ).

1.1. **Definitions.** Section 2.1. “Definitions”, Ch. 2 Definitions and Fundamentals of Biblarz and Sutton (2001) [3]

The **total impulse**  $I_t$  is the thrust force  $F_{\text{thrust}} = F_{\text{thrust}}(t)$  integrated over the burning time  $T$ , i.e.

$$\text{total impulse } I_t = \int_0^T F_{\text{thrust}} dt$$

$t \equiv$  burning time  $\equiv t_p$

specific impulse  $I_s$ , total impulse per unit weight of propellant

$$I_s = \frac{\int_0^t F_{\text{thrust}} dt}{g_0 \int \dot{m} dt} = \frac{I_t}{m_p g_0}$$
$$I_s = \frac{\dot{m} u_e t_p}{\dot{m} g t_p} = \frac{u_e}{g}$$

If  $F_{\text{thrust}} = \dot{m} u_e$ , and constant propellant mass flow,  $u_e$ .

Instead, define

**Definition 1** (total impulse).

(2)

$$I_{tot} \equiv I_t = \int_0^t dt F_{thrust}(t)$$

Let  $m_{p,\text{exp}} = m_{p,\text{exp}}(t)$  = mass of propellant already expelled by time  $t$ .  $m_{p,\text{exp}}(t) \in C^\infty(\mathbb{R})$ .

Biblarz and Sutton (2001) [3] defined *specific impulse* to be the total impulse per unit weight of propellant as such:

**Definition 2** (Specific Impulse).

(3)

$$I_s := \frac{\int_0^t dt F_{thrust}(t)}{W_{p, total \text{ expelled}}}$$

and so

(4)

$$I_s = \frac{\int_0^t dt F_{thrust}(t)}{g_0 \int dt \dot{m}_{p,exp}(t)}$$

For practical calculations,

$$(5) \quad I_s \frac{\overline{F}_{thrust}}{g_0 \overline{\dot{m}}_{p, expelled}}$$

so it's more honest to say that it's the time-averaged specific impulse.

Note that  $\int dt \dot{m}_{p, \exp}(t) = m_{p, \exp} =$  total effective propellant mass expelled through the nozzle.

If we are equipped with the exterior covariant derivative  $D : \Omega^{d-1}(N; TN) \rightarrow \Omega^d(N; TN)$ , then the force on a fluid over any general manifold is

$$F = \int_B \dot{m} \otimes \mathbf{u} + m \otimes \left( \frac{\partial \mathbf{u}}{\partial t} + \nabla_{\mathbf{u}} \mathbf{u} \right)$$

If  $\frac{\partial \mathbf{u}}{\partial t} = 0$ ,  $\nabla_{\mathbf{u}} \mathbf{u} = 0$  (no curved space),

$$F = \int_B \dot{m} \otimes \mathbf{u}$$

If  $\mathbf{u}$  constant over  $B$ , call the *effective exhaust velocity*  $c \equiv u$  to be

$$(6) \quad \begin{aligned} F &= \dot{m} u \\ u &= \frac{F}{\dot{m}} = I_s g_0 \end{aligned}$$

cf. Eq. 2-6 of Biblarz and Sutton (2001) [3]

**Mass ratio**  $\mathbf{MR}$  of a vehicle or *particular vehicle stage* is defined as

$$\mathbf{MR} := \frac{m_f}{m_0}$$

where, using Biblarz and Sutton (2001)'s notation,

$m_f \equiv$  final mass (after rocket operation has consumed all usable propellant), and  $m_0$  (before rocket operation).

$\mathbf{MR}$  applies to a single, or multistage vehicle.

- Overall mass ratio is the product of the individual vehicle stage mass ratios
- Final mass  $m_f$  is mass of vehicle after rocket has ceased to operate when all useful propellant mass  $m_p$  has been consumed and ejected.
- $m_f$  includes all components not useful propellant and may include guidance devices, navigation gear, **payload**

1.1.1. *Mass definitions for rockets, vehicles, propellants.* Instead, let the mass of the rocket + remaining propellant be  $M = M(t) \in C^\infty(\mathbb{R})$ , such that  $M \geq 0$ ,  $\forall t \in \mathbb{R}$ , always.

Let  $M_f \in \mathbb{R}^+$  be a constant, what Biblarz and Sutton calls the final or inert propulsion mass, and denote as  $m_f$ .

Let  $m_p = m_p(t) =$  propellant remaining on vehicle or stage. Thus,

$$(7) \quad \dot{m}_p \leq 0$$

which is expected by physical considerations.

Let

$$(8) \quad m_p(t) = m_p(0) - m_{p, \exp}(t)$$

where  $m_{p, \exp}(t) =$  total propellant that had been expelled. Clearly

$$(9) \quad \dot{m}_{p, \exp}(t) > 0, m_{p, \exp}(0) = 0$$

Thus

$$(10) \quad \begin{aligned} M(t) &= M_f + m_p(t) \\ M(0) &= M_f + m_p(0) \end{aligned}$$

and so define

**Definition 3** (Mass ratio).

$$(11) \quad \mathbf{MR} := \frac{M_f}{M(0)} = \frac{M_f}{M_f + m_p(0)}$$

and

**Definition 4** (propellant mass fraction  $\zeta$ ).

$$(12) \quad \zeta := \frac{m_p(0)}{M(0)} = \frac{M(0) - M_f}{M(0)} = \frac{m_p(0)}{M_f + m_p(0)}$$

Biblarz and Sutton (2001) uses a different notation for the masses.

e.g.  $\zeta = \frac{m_p(0)}{M(0)} = 0.91$  means only 9 percent of mass is inert rocket hardware, and this small fraction contains, feeds, burns a substantially larger mass of propellant: a high value of  $\zeta$  is desirable.

Impulse to weight ratio of a complete propulsion system:

$$\frac{I_t}{w_0}$$

$I_t \equiv$  total impulse,

$w_0 \equiv$  initial or propellant-loaded vehicle weight.

It's reasonable to suppose that

$$\frac{I_t}{w_0} = \frac{I_t}{M(0)g} = \frac{I_t}{(M_f + m_p(0))g_0} = \frac{I_s m_p(0)g_0}{(M_f + m_p(0))g_0} = \frac{I_s}{\frac{M_s}{m_p} + 1}$$

thrust-to-weight ratio  $F/w_0$

Example 2-1 of Biblarz and Sutton (2001) [3] has some subtleties about the mass definitions that wasn't covered clearly (at least for me) by the discussion.

Let  $M_r$  is the mass of the propulsion system,  $r$  standing for "rocket." If there is *no payload*, then

$$M(0) = M_r + m_p(0)$$

$$M(t) = M_r + m_p(t)$$

*But*, if payload is considered, observe that payload can be considered to be part of the whole of the vehicle (imagine an imaginary box surrounding the entire system or "vehicle"), *and also* can be considered to be separate and distinct (no overlap) from the propulsion system plus propellant (imagine an imaginary box surrounding only the rocket (propulsion system) and fuel, and another imaginary box surrounding the payload; the two imaginary boxes do not overlap).

$$M(0) = M_r + M_{\text{pay}} + m_p(0)$$

$$M(t) = M_r + M_{\text{pay}} + m_p(t)$$

Recall that mass ratio is defined as  $\mathbf{MR} := \frac{M_f}{M(0)}$ . But remember that  $\mathbf{MR}$  can be considered for different, discrete imaginary boxes (subsystems). So for Example 2-1 of Biblarz and Sutton:

$$\mathbf{MR}_{\text{tot}} = \frac{M_{f, \text{tot}}}{M(0)} = \frac{M_r + M_{\text{pay}}}{M(0)}$$

$$\mathbf{MR}_r = \frac{M_r}{M_r + m_p(0)} = \frac{M_r}{M(0) - M_{\text{pay}}}$$

Notice how you can not take the payload mass into account; imagine only dealing with what's in the imaginary box surrounding only the rocket plus propellant.

If  $\frac{\partial \mathbf{u}}{\partial t} = 0$ , assuming flat space,  $F = \int_S \dot{m}_S \otimes \mathbf{u}$ . Notice that for the units to make sense,  $\dot{m}_S$  is the mass per unit (surface) area. If  $\dot{m}_S, \mathbf{u}$  constant over surface area  $S$ ,  $F = \dot{m} \mathbf{u}$ . Then,  $u = \frac{F}{\dot{m}} = I_s g_0$ .

Then  $\bar{u} = I_s g_0$ .  $I_{\text{tot}} = I_s m_p(0) = I_s (M(0) - (M_r + M_{\text{pay}}))$

Given the burn duration  $T$ , estimate  $\dot{m}_{p,\text{exp}}$  with the simplest (linear) approximation:  $\dot{m}_{p,\text{exp}} = \frac{m_p(0)}{T}$ .

$$\bar{F}_{\text{thrust}} = \bar{\dot{m}}_{p,\text{exp}} g_0 I_s$$

Consider  $\bar{F}_{\text{thrust}}$  on rocket, only. For a horizontal trajectory, max. acceleration is found at end of the thrusting schedule, just before shutdown (because while thrust is unchanged the mass is now at its minimum value). Think of "Max-Q".

$$\implies a_f = \frac{\bar{F}_{\text{thrust}}}{M_f}$$

**1.2. Definitions with Thrust; effective exhaust velocity revisited.** Thrust is the force produced by a rocket propulsion system acting upon a vehicle.

In a simplified way, it's the reaction experienced by the structure due to ejection of matter at high velocity.

It represents the same phenomenon that pushes garden hose backwards or makes a gun recoil.

In the latter case, forward momentum of the bullet and powder charge is equal to recoil or rearward momentum of gun barrel.

From pp. 30, Sec. 2.1. "Definitions" of Biblarz and Sutton (2001) [3], effective exhaust velocity  $c$  was defined:

**Definition 5** (effective exhaust velocity). *effective exhaust velocity  $c :=$  average equivalent velocity at which propellant is ejected from vehicle.*

Further, Biblarz and Sutton (2001) makes the assumption of uniform axial velocity  $c$ . But, in a rocket nozzle, the actual exhaust velocity is not uniform over the entire exit cross section, and doesn't represent the entire thrust magnitude. Nevertheless, Biblarz and Sutton (2001) gives the following formulas:

$$c = I_s g_0 = F / \dot{m}$$

cf. Eq. 2-6 of Biblarz and Sutton (2001) [3].

Let's derive the thrust, due to the change in momentum.

**1.2.1. Momentum flux as a differential form.** Consider this object:

$$u_x dx \wedge dy \wedge dz + u_y dy \wedge dz \wedge dx + u_z dz \wedge dx \wedge dy$$

It is invariant under transformation  $P$  such that  $(x, y, z) \xrightarrow{P} (-x, -y, -z)$ .

It is also invariant under  $x \mapsto -x$  **and**  $\mathbf{u} \mapsto -\mathbf{u}$ , where  $x$  could be  $x, y$  or  $z$ .

Now

$$u_x dx \wedge dy \wedge dz + u_y dy \wedge dz \wedge dx + u_z dz \wedge dx \wedge dy = \sum_{i=1}^3 \left( u_i dx^i \wedge \frac{\epsilon_{ijk}}{2!} dx^j \wedge dx^k \right) = \sum_{i=1}^3 \frac{u_i \epsilon_{ijk}}{2!} dx^i \wedge dx^j \wedge dx^k$$

While the last expression is "manifestly covariant" in the indices  $j$  and  $k$  (i.e. the pair of "up" and "down" indices "match up"), it doesn't appear the case for the  $i$  index.

**1.2.2. Mass conservation.** Recall the derivation of mass conservation:

For time  $t \in \mathbb{R}$  being a 1-dim. parameter,

given mass density  $\rho = \rho(t, x) \in \mathbb{R} \times N$ ,  $x \in N$ ,

volume form  $\text{vol}^n \in \Omega^n(N)$

on a smooth (spatial) manifold  $N$ , smooth submanifold with boundary  $B(t) \subset N$ , representing the control volume,

total mass inside a control volume  $B(t)$ ,  $m = m(t)$ , is given by

$$m := \int_{B(t)} \rho \text{vol}^n$$

Then for mass conservation,

$$\begin{aligned} \dot{m} &\equiv \frac{d}{dt} m = \frac{d}{dt} \int_{B(t)} \rho \text{vol}^n = \int_{B(t)} \mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}}(\rho \text{vol}^n) = \int_{B(t)} \left[ \frac{\partial \rho}{\partial t} \text{vol}^n + \mathcal{L}_{\mathbf{u}}(\rho \text{vol}^n) \right] = \int_{B(t)} \left[ \frac{\partial \rho}{\partial t} \text{vol}^n + \mathbf{d}(i_{\mathbf{u}} \rho \text{vol}^n) \right] = \\ (13) \quad &= \int_{B(t)} \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial B} i_{\mathbf{u}} \rho \text{vol}^n = \\ &= \int_{B(t)} \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial B} \rho u^i dS_i \end{aligned}$$

where

$$\mathcal{L}_{\mathbf{u}}(\rho \text{vol}^n) = (\mathbf{d}i_{\mathbf{u}} + i_{\mathbf{u}}\mathbf{d})(\rho \text{vol}^n) = \mathbf{d}(i_{\mathbf{u}}\rho \text{vol}^n) + 0 = \mathbf{d}i_{\mathbf{u}}\rho \text{vol}^n$$

is due to Cartan's magic formula, and

$$\int_{B(t)} \mathbf{d}(i_{\mathbf{u}}\rho \text{vol}^n) = \int_{\partial B(t)} i_{\mathbf{u}}\rho \text{vol}^n$$

is by Stoke's law.

If  $\frac{\partial \rho}{\partial t} = 0$  (steady-state condition), then

$$\dot{m} = \int_{\partial B} \rho u^i dS_i$$

## 2. NOZZLE THEORY AND THERMODYNAMIC RELATIONS

**2.1. Isentropic Flow Through Nozzles.** Subsection 3.3 "Isentropic Flow Through Nozzles" of Biblarz and Sutton (2001) [3]

The Bernoulli invariant for compressible flow <sup>1</sup> gives us this:

$$\begin{aligned} h_1 + \frac{1}{2}v_1^2 &= h_2 + \frac{1}{2}v_2^2 \text{ or } v_2 = \sqrt{2(h_1 - h_2) + v_1^2} \\ \implies h_1 - h_2 &= \frac{C_p}{MN}(\tau_1 - \tau_2) = \frac{\gamma}{\gamma - 1} \frac{\tau_1}{M} \left( 1 - \frac{\tau_2}{\tau_1} \right) = \frac{\gamma}{\gamma - 1} \frac{\tau_1}{M} \left( 1 - \left( \frac{p_2}{p_1} \right)^{\frac{\gamma-1}{\gamma}} \right) = \frac{\gamma R T_1}{\gamma - 1} \left( 1 - \left( \frac{p_2}{p_1} \right)^{\frac{\gamma-1}{\gamma}} \right) \end{aligned}$$

so that

$$v_2 = \sqrt{\frac{2\gamma R T_1}{\gamma - 1} \left( 1 - \left( \frac{p_2}{p_1} \right)^{\frac{\gamma-1}{\gamma}} \right) + v_1^2}$$

When chamber section is large compared to nozzle throat section, chamber velocity or nozzle approach velocity comparatively small and  $v_1^2$  can be neglected.[3]

Chamber temperature  $T_1$  at nozzle inlet; under isentropic conditions,  $T_1$  differs little from stagnation temperature or (for chemical rocket) combustion temperature.[3]

Ch. 4 Flight Performance of Biblarz and Sutton (2001) [3] should probably be read before Ch. 3 Nozzle Theory and Thermodynamic Relations or before. In AE121 Fall 2015, the material for Ch. 4 was covered in lectures (I think, from the problem sets) before Thermodynamics and Nozzle Theory.

[3]

## 3. FLIGHT PERFORMANCE

cf. Chapter 4 Flight Performance of Biblarz and Sutton (2001) [3]

<sup>1</sup> "Euler equations (fluid dynamics)", Wikipedia

### 3.1. Gravity-Free, Drag-Free Space Flight.

$m_p \equiv$  (total) propellant mass (initially)

$t_p \equiv$  propellant burning duration time

$$m(0) = 0$$

Let mass of rocket + propellant  $M = M(t) = M(0) - m(t)$  s.t.  $M(0) = M_0 + m_p$

$$m(t_p) = m_p$$

$$M(t_p) = M_0$$

If  $\dot{m} = \frac{m_p}{t_p}$  (assume constant propellant flow rate),

$$M = M(0) - \frac{m_p}{t_p}t = M(0) \left(1 - \frac{m_p}{M(0)} \frac{t}{t_p}\right) = M(0) \left(1 - \left(1 - \frac{M(0) - m_p}{M(0)}\right) \frac{t}{t_p}\right)$$

cf. Eq. 2-7 of Biblarz and Sutton (2001) [3].

Define some quantities: mass ratio  $= \frac{m_f}{m_0}$  where

$m_f :=$  final mass (after rocket operation had consumed all usable propellant), which is  $M_0$  above

$m_0 :=$  mass before rocket operation, which is  $M(0)$

propellant mass fraction  $\frac{m_p}{M(0)}$  cf. Eq. 2-8 of Biblarz and Sutton (2001) [3].

For thrust  $F_{\text{thrust}}$ ,

$$\begin{aligned} F_{\text{thrust}} = \dot{m}u_e &= M \frac{du}{dt} \implies \frac{\dot{m}dt}{M} = \frac{du}{u_e} \text{ or } \frac{\Delta u}{u_e} = \int \frac{\frac{m_p}{t_p}dt}{M(0) - \frac{m_p}{t_p}t} = -\ln(M(0) - \frac{m_p}{t_p}t) = \\ &= -\ln(M(0) - \frac{m_p}{t_p}t_p) + \ln(M(0)) = \ln\left(\frac{M(0)}{M(0) - \frac{m_p}{t_p}t_p}\right) = \ln\left(\frac{M(0)}{M(0) - m_p}\right) \end{aligned}$$

Thus

$$\begin{aligned} \Delta u &= u_e \ln\left(\frac{M(0)}{M(0) - m_p}\right) = u_e \ln\left(\frac{M(0)}{M_0}\right) \text{ or} \\ \exp\left(\frac{\Delta u}{u_e}\right) &= \frac{M(0)}{M(0) - m_p} = \frac{M(0)}{M_0} \text{ or } \frac{M_0}{M(0)} = \exp\left(\frac{-\Delta u}{u_e}\right) \end{aligned}$$

Also remember that  $F = \dot{m}u_e = \dot{m}I_{\text{sp}}g_0$ , so  $u_e \equiv$  effective exhaust velocity can be related directly to the specific impulse  $I_{\text{sp}}$ .

Also note that for propellant mass fraction  $\frac{m_p}{M(0)} < 1$ ,  $\frac{m_p}{M(0)} = \frac{M(0) - M_0}{M(0)} = 1 - \frac{M_0}{M(0)}$

### 3.2. Forces Acting on a Vehicle in the Atmosphere.

assume starting and stopping transients very short and neglected,  $F = u_e \dot{m}$

if mass rate of propellant consumption  $\dot{m}$  constant,  $\dot{m} = \frac{m_p}{t_p}$  so  $F = \frac{m_p}{t_p}u_e$

drag  $D$  opposite to flight path due to resistance of body to motion in fluid

lift  $L$  normal to flight path

$$L = C_L \frac{1}{2} \rho A u^2$$

$$D = C_D \frac{1}{2} \rho A u^2$$

cf. Eqns. (4-10), (4-11) of Biblarz and Sutton (2001) [3]

density of earth's atmosphere can vary by factor up to 2 (for altitudes of 300 to 1200 km)

depending on solar activity and night-to-day temperature variations. major unknown in drag

Assume neglect variation of gravity with geographical features and oblate shape of earth,

$$\begin{aligned} F &= \frac{GM_e m}{r^2} \\ mg_0 &= \frac{GM_e}{R_0^2} m \implies g = g_0 \frac{R_0^2}{R^2} = g_0 \left(\frac{R_0}{R_0 + h}\right)^2 \end{aligned}$$

### 3.3. Basic Relations of Motion.

cf. Section 4.3 “Basic Relations of Motion,” Chapter 4 Flight Performance of Biblarz and Sutton (2001) [3]

$$m\dot{\mathbf{u}} = \mathbf{F} - \mathbf{D} - m\mathbf{g}$$

$$m\dot{u} = F \cos(\psi - \theta) - D - mg \sin \theta$$

$$\implies mu \frac{d\theta}{dt} = F \sin(\psi - \theta) + L - mg \cos \theta$$

$\psi =$  direction of thrust angle from horizontal reference

Consider perturbation effects (cf. Sec. 4.6, listed), drag and gravity

$$\begin{aligned} \dot{u} &= \frac{F}{m} \cos(\psi - \theta) - \frac{C_D}{2m} \rho A u^2 - g \sin \theta \\ \implies u\dot{\theta} &= \frac{F}{m} \sin(\psi - \theta) + \frac{C_L}{2m} \rho A u^2 - g \cos \theta \end{aligned}$$

$C_D, C_L$  are functions of velocity or Mach number (!!!)

For actual trajectory analyses, perturbation effects in Sec. 4.6 must be considered

3-body theory considered

when propellant flow and thrust not constant

from optical or radar tracking data, thrust or actual specific impulse during actual vehicle flights determined from accurately observed trajectory data.

make assumption or measurement on propellant flow (which usually varies in a predetermined manner)

If  $L = 0$  (wingless rocket projectile),  $\psi - \theta = 0$  (flight direction  $\theta$  same as thrust direction), and

$M(t) = M(0) - \dot{m}t$ ; assume constant  $\dot{m} = \frac{t}{t_p}m_p$

$M(t) = M(0) - \frac{tm_p}{t_p} = M(0)(1 - \frac{m_p}{M(0)} \frac{t}{t_p})$ , with  $\xi \equiv \frac{m_p}{M(0)}$  propellant mass ratio

$$\begin{aligned} \dot{u} &= \frac{F}{M} - \frac{C_D}{2M} \rho A u^2 - g \sin \theta = \frac{u_e \dot{m}}{M} - \frac{C_D \rho A u^2}{2M} - g \sin \theta = \\ &= \frac{u_e \xi \frac{t}{t_p}}{1 - \xi \frac{t}{t_p}} - \frac{C_D \rho A u^2}{2M(0)(1 - \xi \frac{t}{t_p})} - g \sin \theta \end{aligned}$$

and

$$u\dot{\theta} = -g \cos \theta$$

**Example 4-1.** Consider “a simple-stage rocket for a rescue flare has the following characteristics and its flight path nomenclature is shown in the sketch.”

Neglect drag “since flight velocities are low,” assume no wind, assume local acceleration of gravity to be equal to sea level  $g_0$ , invariant throughout flight.

**Problems. Problem 1.** Recall that

$$\frac{M_0}{M(0)} = \exp\left(\frac{-\Delta u}{u_e}\right)$$

and so (in Python)



```
import sympy
from sympy import *
>>> exp(-1600/2000.)
0.449328964117222
```

Problem 2.

$$\frac{m_p}{M(0)} = \frac{M(0) - M_0}{M(0)} = 1 - \frac{M_0}{M(0)} = 1 - \frac{1}{5} = 4/5 = 0.8$$

Problem 3. dragless projectile, so  $D = 0$ .

$$\dot{u} = -g_0 + \frac{u_e}{t_p(M(0)/m_p - t/t_p)} \implies \Delta u = -g_0 t - u_e \ln(1 - \frac{m_0 t}{M(0)t_p})$$

Plugging in  $u_e = 2209\,m/\text{sec}$ ,  $m_p/M(0) = 0.57$ ,  $t_p = 5.0\,\text{sec}$ ,  $u_0 = h_0 = 0$ ,

```
Propulsion.py
>>> integrate( flightpathdirection.subs(Drag,0).subs(psi,theta).subs(theta,pi/2).subs(F_thrust, m_p/t_p*u_e).
subs(M,M_constantflow).subs(m_p,M0*0.57).subs(t_p,5.).subs(u_e,2209.).subs(g_0,9.8).factor(M0).rhs, (t,0,5.0) )
1815.32988528061
```

```
Problem0403 = flightpathdirection.subs(Drag,0).subs(psi,theta).subs(theta,pi/2).subs(F_thrust, m_p/t_p*u_e).
subs(M,M_constantflow).subs(m_p,M0*0.57).subs(t_p,5.).subs(u_e,2209.).subs(g_0,9.8).factor(M0).rhs
```

$u_p = 1815\,m/\text{sec}$

```
>>> integrate( integrate(Problem0403,(t,0,t) ),(t,0,5.0) )
3890.37850288891
```

$h_p = 3.89 \times 10^3$

Problem 4. How to estimate  $A$  of the projectile?

Problem 5.

Now

$$M(t)a = F_{\text{thrust}}$$
$$I_{sp} = \frac{F_{\text{thrust}}t_p}{g_0 \frac{m_p}{t_p}} = \frac{F_{\text{thrust}}t_p}{g_0 m_p}$$

and for  $M(t) = M(0) - \frac{m_p}{t_p}t$ ,

$$a = \frac{I_{sp} \left( \frac{g_0 m_p}{t_p} \right)}{M(0) - \frac{m_p}{t_p}t} \leq a(t = t_p) \implies a(t = t_p) = \frac{I_{sp} \left( \frac{g_0 m_p}{t_p} \right)}{M_0} = \frac{I_{sp} g_0}{t_p} \left( \frac{1}{\frac{M(0)}{m_p} - 1} \right)$$

Plugging  $a = 50\,m/\text{sec}^2$  and solving for  $t_p$ ,

```
>>> 260.*9.8/50.*(1/ ( 1/0.88 - 1 ) )
373.70666666666637
```

(a)

$373.71\,\text{sec}$

maximum allowable burn time, assuming steady propellant mass flow

(b)

$$a = \frac{I_{sp} g_0 \xi / t_p}{1 - \xi t / t_p} \implies \Delta u = -I_{sp} g_0 \ln(1 - \xi t / t_p)$$

so

$$\Delta u = 5402.4\,m/\text{sec}$$

for maximum velocity relative to the launch vehicle

Problem 6. Satellite in circular orbit

$$F = \frac{GM_e m}{(R_0 + h)^2} = \frac{mv^2}{(R_0 + h)} \implies \sqrt{\frac{GM_0}{R_0 + h}} = v$$
$$\frac{2\pi(R_0 + h)}{T} = v \implies T = \frac{2\pi(R_0 + h)^{3/2}}{\sqrt{GM_e}}$$
$$\frac{1}{2}mv^2 + \frac{-GM_e m}{R_0 + h} - \left( \frac{-GM_e m}{R_0} \right) = m \left[ \frac{1}{2} \frac{GM_e}{R_0 + h} - \frac{GM_e}{R_0 + h} + \frac{GM_e}{R_0} \right] = m \left[ GM_e \left( \frac{-1}{2(R_0 + h)} + \frac{1}{R_0} \right) \right]$$

Then run **Propulsion.py** which now imports (**import**) in Physique, a small package with the NIST (National Institute of Standards and Technology) Fundamental Constants **FundConst**, NIST SI conversions **conv**, and NASA Planetary Fact Sheet **plnfacts** as Python **pandas** DataFrames.

```
M_earth = plnfacts.loc[plnfacts['Planet']=="EARTH", "Mass_(1024kg)"].values[0]*10**(24) # in kg
R_earth = plnfacts.loc[plnfacts['Planet']=="EARTH", "Diameter_(km)"].values[0]/Decimal(2)
```

```
Gconst = FundConst[ FundConst["Quantity"].str.contains("gravitation") ].loc[243,"Value"]
v0406 = sqrt( Gconst*M_earth/((R_earth + Decimal(500))*10**3) )
# velocity of satellite v of Chapter 4, Problem 6 of Biblarz and Sutton
7611.17633707692
```

```
T0406 = (2.*N(pi)*float(((R_earth + Decimal(500))*10**3 )**(3./2))/float(sqrt( Gconst*M_earth))
# 5678 secs. or 1.58 hours
```

```
Eperm0406 = Gconst*M_earth*(-1/(2*((R_earth+Decimal(500))*10**3)) + 1/(R_earth*10**3))
# Energy per mass
'%6E' % Eperm0406 # 33.51 MJ/kg
```

$v = 7611\,m/\text{sec} \quad T = 5678\,s \text{ or } 1.58\,\text{hours} \quad 33.51\,MJ/kg$

## Part 2. Notes and Solutions for *Space Propulsion Analysis and Design* by Humble, Henry, Larson

cf. Humble, Henry, and Larson (1995) [4]

### 4. THERMODYNAMICS OF FLUID FLOW; CF. CH. 3. OF HUMBLE, HENRY, AND LARSON (1995) [4]

cf. Humble, Henry, and Larson (1995), Ch. 3. Thermodynamics of Fluid Flow [4]

Chamber

Chamber conditions

$$(p_c, T_c, \rho_c) \in \mathbb{R}^2 \times \mathbb{R}_+$$

$p_c$  = chamber pressure

$T_c$  = chamber temperature

$\rho_c$  = chamber density

Chamber characteristics

- Combustion
- high pressure  $p_c$
- high temperature  $T_c$
- very low net fluid velocity

Nozzle.

exit conditions  $(p_e, T_e, \rho_e) \in \mathbb{R}^2 \times \mathbb{R}_+$

$p_e$  = exit pressure

$T_e$  = exit temperature

$\rho_e$  = exit density

Exit characteristics.

- flow expands to fill enlarged volume
- reduced  $p_e$
- reduced  $T_e$
- very high fluid velocity

Generalized View of a Rocket Thrust Chamber:  
Understanding the thermodynamic conditions is key to understanding its performance

4.1. **3.1 Mass Transfer.** mass flow rate perpendicular to duct of cross-sectional area (A) is

$$\dot{m} = \rho v A$$

$\dot{m}$  = mass flow rate (kg/s)  
 $\rho$  = fluid density (kg/m<sup>3</sup>)  
 $v$  = fluid velocity ( $m/s$ )  
 $A$  = cross-sectional area of duct (m<sup>2</sup>)

Control volume - region with constant shape and size that stays fixed in space  
system approach (a.k.a. control mass approach) focus on fixed amount of matter  
envelope containing the matter may change its size and shape, location  
mass conservation (control-volume approach)

$$\frac{d}{dt}(m_{cv}) = \dot{m}_{\text{in}} - \dot{m}_{\text{out}}$$

$m_{cv}$  = amount of mass in control volume (kg)  
 $\dot{m}_{\text{in}}$  = mass flow rate into control volume (kg/s)  
 $\dot{m}_{\text{out}}$  = mass flow rate out of control volume (kg/s)

$$\dot{m} = \int_{cs} \rho v dA$$

$cs$  = control surface  
 $\dot{m}$  = total mass flow through control surface (kg/s)  
 $\rho$  = fluid density in flow (kg/m<sup>3</sup>)  
 $v$  = velocity of particles in flow (m/s)  
 $dA$  = elemental area in control surface ( $m^2$ )

4.2. **3.2 Thermodynamic Relations (Energy and Entropy).** cf. Humble, Henry, and Larson (1995), Ch. 3. Thermodynamics of Fluid Flow [4]

4.2.1. *3.2.5. Isentropic Flow in One-Dimension.* cf. Humble, Henry, and Larson (1995), Ch. 3. Thermodynamics of Fluid Flow [4], pp. 95

**Simplifying Assumptions**

Fluid flow in a duct, applies to gaseous propellant in a nozzle, liquid or gaseous propellant in fuel plumbing, and coolant through nozzle’s cooling tubes

- (1) isentropic flow, i.e. reversible and adiabatic; from the chamber (after combustion of propellant) to the nozzle exit.
  - *adiabatic flow* means heat transfer doesn’t dissipate energy from flow
  - although heat transfer through chamber walls is significant, it’s relatively small percentage of the total energy generated.
  - This assumption also neglects effects of friction and fluid viscosity and doesn’t apply to shock waves
- (2) 1-dim. flow.

larger nozzle cone half angle, less acceptable this assumption becomes. For most nozzles, assumption causes less than

5

- (3) products of combustion constitute a perfect gas

- (4) Frozen flow. Once established in chamber, products of combustion don’t change in chemical composition while traversing nozzle

- (5) steady flow

*Isentropic Relations*

energy conservation  $\implies Q = dU - W = dU + pdV$

reversible process  $\implies Q = \tau d\sigma$

$$\tau d\sigma = dU + pdV$$

$$dU = C_V d\tau$$

$$d\sigma = \frac{C_V}{\tau} d\tau + \frac{p}{\tau} dV = \frac{C_V}{\tau} d\tau + \frac{N}{V} dV$$

$d\sigma = 0$  isentropic process

$$\implies 0 = \frac{C_V}{\tau} d\tau + \frac{N}{V} dV \text{ or } \frac{C_V}{\tau} d\tau = -\frac{N}{V} dV$$

$$\begin{aligned} C_V \ln(\tau_2/\tau_1) &= N \ln(V_1/V_2) = \\ &= (C_p - C_V) \ln(V_1/V_2) \\ C_p &= C_V + N \implies \\ \implies \ln(\tau_2/\tau_1) &= (\gamma - 1) \ln(V_1/V_2) \\ \implies \frac{\tau_2}{\tau_1} &= \left(\frac{V_1}{V_2}\right)^{\gamma-1} = \left(\frac{\rho_2}{\rho_1}\right)^{\gamma-1} \end{aligned}$$

with

$$\begin{aligned} \gamma &\equiv \frac{C_p}{C_v} \\ \gamma - 1 &= \frac{N}{C_V} \end{aligned}$$

And so for *isentropic, 1-dim., steady flow* and calorically perfect gas:

$$(14) \quad \begin{aligned} (\tau_1, \rho_1) &\mapsto (\tau_2, \rho_2) \\ (\tau_1, p_1) &\mapsto (\tau_2, p_2) \\ (\rho_1, p_1) &\mapsto (\rho_2, p_2) \end{aligned} \quad \text{with} \quad \begin{aligned} \frac{\tau_2}{\tau_1} &= \left(\frac{\rho_2}{\rho_1}\right)^{\gamma-1} \text{ or } \frac{\rho_1^{\gamma-1}}{\tau_1} = \frac{\rho_2^{\gamma-1}}{\tau_2} \\ \frac{p_2}{p_1} &= \left(\frac{\tau_2}{\tau_1}\right)^{\frac{\gamma}{\gamma-1}} \\ \frac{p_2}{p_1} &= \left(\frac{\rho_2}{\rho_1}\right)^{\gamma} \end{aligned}$$

From the Bernoulli invariant, and assuming

- (1) steady flow
- (2) adiabatic process (no heat transfer)
- (3) no significant changes in potential energy
- (4) No shaft work or shear work done

$$(15) \quad h_1 + \frac{1}{2}u_1^2 = h_2 + \frac{1}{2}u_2^2$$

From

$$\begin{aligned} \mathfrak{M} &:= \frac{u}{a} \\ a^2 &= \frac{\gamma \tau}{M} \quad \text{and} \quad C_p = \frac{N}{\gamma - 1} + N = N \left( \frac{\gamma}{\gamma - 1} \right) \implies \\ c_p &\equiv \frac{C_p}{MN} \quad \frac{C_p}{N} = \frac{\gamma}{\gamma - 1} \\ c_p(\tau_1 - \tau_2) &= \frac{1}{2}(u_2^2 - u_1^2) = \\ \implies &= \frac{1}{2}(\mathfrak{M}_2^2 \gamma \tau_2 - \mathfrak{M}_1^2 \gamma \tau_1) \frac{1}{M} \\ \implies c_p \tau_1 + \frac{1}{2} \mathfrak{M}_1^2 \gamma \frac{\tau_1}{M} &= (c_p + \frac{1}{2} \mathfrak{M}_2^2 \frac{\gamma}{M}) \tau_2 \\ \implies \frac{\tau_1}{\tau_2} &= \frac{\frac{C_p}{N} + \frac{1}{2} \mathfrak{M}_2^2 \gamma}{\frac{C_p}{N} + \frac{1}{2} \mathfrak{M}_1^2 \gamma} = \frac{\frac{\gamma}{\gamma-1} + \frac{1}{2} \mathfrak{M}_2^2 \gamma}{\frac{\gamma}{\gamma-1} + \frac{1}{2} \mathfrak{M}_1^2 \gamma} = \frac{1 + \frac{1}{2} \mathfrak{M}_2^2 (\gamma - 1)}{1 + \frac{1}{2} \mathfrak{M}_1^2 (\gamma - 1)} \end{aligned}$$

#### Area ratio

From Eq. 35,

**Definition 6** (Expansion ratio).

$$(16) \quad \epsilon \equiv \text{expansion ratio} := \frac{A_{exh}}{A_t} = \frac{1}{\mathfrak{M}_{exh}} \left[ \left( 1 + \frac{\gamma - 1}{2} \mathfrak{M}_{exh}^2 \right) \frac{2}{\gamma + 1} \right]^{\frac{\gamma+1}{2(\gamma-1)}}$$

So from

$$\text{mass conservation} \implies \dot{m} = \rho_t u_t A_T$$

$$\text{Mach number (definition)} \implies \mathfrak{M}_t := u_t / a_t$$

$$\begin{aligned} \frac{\rho_t \tau_0}{M} &= \frac{\rho_0 \left( \frac{2}{\gamma+1} \right)^{\frac{1}{\gamma-1}} \tau_0}{M} = p_0 \left( \frac{2}{\gamma+1} \right)^{\frac{1}{\gamma-1}} \\ \frac{1}{\gamma-1} + \frac{1}{2} &= \frac{2+\gamma-1}{\gamma-1} = \frac{\gamma+1}{\gamma-1} \end{aligned}$$

$$\begin{aligned} \dot{m} u_{\text{exh}} &= \rho_t \mathfrak{M}_t a_t A_t u_{\text{exh}} = \rho_t \left( \frac{\gamma \tau_t}{M} \right)^{1/2} A_t \left[ \frac{2\gamma}{\gamma-1} \frac{\tau_0}{M} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right) \right]^{1/2} = \frac{\rho_t A_t}{M} \left[ \frac{2\gamma^2 \tau_t \tau_0}{\gamma-1} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right) \right]^{1/2} = \\ &= p_0 \left( \frac{2}{\gamma+1} \right)^{\frac{1}{\gamma-1}} A_t \left[ \frac{2\gamma^2}{\gamma-1} \frac{2}{\gamma+1} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right) \right]^{1/2} = A_t p_0 \left[ \frac{2\gamma^2}{\gamma-1} \left( \frac{2}{\gamma+1} \right)^{\frac{\gamma+1}{\gamma-1}} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right) \right]^{1/2} \end{aligned}$$

So from

$$F_{\text{thrust}} = \dot{m} u_{\text{exh}} + (p_e - p_a) A_e$$

$$(17) \quad \boxed{F_{\text{thrust}} = A_t p_0 \left[ \frac{2\gamma^2}{\gamma-1} \left( \frac{2}{\gamma+1} \right)^{\frac{\gamma+1}{\gamma-1}} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right) \right]^{1/2} + (p_e - p_a) A_{\text{exh}}}$$

Eq. 17 is Eq. (3.129) of Humble, Henry, and Larson (1995), Ch. 3. Thermodynamics of Fluid Flow [4], pp. 112, Eq. (3-29) of Biblarz and Sutton (2001) [3], Ch. 3, pp. 63.

### Part 3. 1-dim. propulsion (revisited)

#### 5. FLUID FLOW (REVIEW)

**5.1. Mass conservation, revisited.** Let  $m$  also denote the  $d = 3$  (differential) form,  $m \in \Omega^d(N)$ , where  $N \equiv$  spatial manifold,  $\dim N = d = 3$ . Let  $V$  be the control volume. Then

$$M = \int_V m = \int_V \rho \text{vol}^d$$

Recall then Eq. 13, so that

$$\dot{M} \equiv \frac{d}{dt} M = \frac{d}{dt} \int_V \rho \text{vol}^d = \int_V \frac{\partial \rho}{\partial t} \text{vol}^d + \int_{\partial V} \rho u^i dS_i$$

**5.2. Newton's second law on continuum, fluid flow.** For  $(\mathbf{d}, -) : \Omega^{d-1}(N) \otimes TN \rightarrow \Omega^d(N) \otimes TN$ , where

smooth (spatial) manifold  $N$ ,

$\Omega^d(N)$  is the space of all  $d$ -forms on  $N$ ,

$TN \equiv$  tangent space on  $N$ ,

$\Omega^d(N) \otimes TN$  is the tensor product space of tensor products of  $d$  differential forms and tangent vectors, then

$$F = \frac{d}{dt} \Pi \equiv \dot{\Pi} := \frac{d}{dt} \int_V m \otimes u = \int_V \dot{m} \otimes u + m \left( \frac{\partial u}{\partial t} + u^k \frac{\partial u}{\partial x^k} \right) = \int_V \frac{\partial(\partial u)}{\partial t} \text{vol}^d + \int_{\partial V} \rho u^k dS_k \otimes u$$

where, recall,

$$\frac{\partial(\rho u)}{\partial t} \text{vol}^d = \frac{\partial \rho}{\partial t} \text{vol}^d \otimes u + \rho \text{vol}^d \otimes \frac{\partial u}{\partial t} = \dot{m} \otimes u + m \frac{\partial u}{\partial t}$$

Hill and Peterson (1992) [2] says it's important to note that  $F \equiv \sum F$  is the sum of the forces applied to the control volume by its environment.

$\int_V \frac{\partial(\rho u)}{\partial t} \text{vol}^d$  is interpreted as the changing of the total momentum instantaneously contained within the control volume  $V$ .

$\int_{\partial V} \rho u^k dS_k \otimes u$  is interpreted as the changing of the net flow rate of momentum leaving the control volume.

Practical advice from "Newton's Second Law and the Momentum Equation", pp. 26, Hill and Peterson (1992) [2]:

- Steady flow means velocity, density, and "so on" (what does Hill and Peterson (1992) [2] mean by "so on") at any point in space don't change with time (though they may well vary from point to point in space). Thus for steady flow  $\frac{\partial(\rho u)}{\partial t} = 0$ , so one can do *complete thermodynamic analysis* by considering only inputs and outputs of a control volume, not its contents.
- If one can convert an unsteady flow problem to a steady one by moving coordinate system of observer, it's usually a good idea.

#### 6. IDEAL ROCKET EQUATION

Sec. 1, Humble, Henry, and Larson (1995)

Looking at pp. 40-41 of [1], Chapter 2, "Thermodynamics and Fluid Flows",

$$(\rho u^2 A)_2 - (\rho u^2 A)_1 = \int_{A_2 \amalg A_1 \amalg \text{side}} T^{ij} dS_j \otimes e_1$$

$\int_{A_2 \amalg A_1 \amalg \text{side}} T^{ij} dS_j \otimes e_1$  comprise of the contribution from the fluid being a perfect fluid,  $T = -pg$ , so consider that *only* first:

$$\int_{A_2 \amalg A_1} -pg^{1j} dS_j \otimes e_1 = -pA_2 + pA_1$$



Suppose shear stress (which happens at solid boundary; by no slip condition, there’s a boundary layer), occurs at wall.

$$\int_{\text{side}} T^{1j} dS_j \otimes e_1 = \int_{\text{side}} (T^{12}(2dx^3 \wedge dx^1) + T^{13}(2dx^1 \wedge dx^2))e_1$$

Assume average (constant) shear over “infinitesimal” or differential cross section:  $T^{12} = T^{13} = \frac{\tau}{2}$

$$\implies \tau \int_{\text{side}} dx^3 \wedge dx^1 + dx^1 \wedge dx^2 = \tau(\text{circumference}) \int dx^1$$

Thus, this explains the last equation on pp. 40, and first 2 equations, including Eq. (2.62) on pp.41 of “Thermodynamics and Fluid Flows”, Chapter 2, of Gas Turbine and Rocket Propulsion of Oates (1997) [1]

Oates (1997) [1]

$$\begin{aligned} \tau d\sigma &= dU + pdV \\ \implies \tau d\sigma &= \left( \frac{\partial U}{\partial \tau} \right)_V d\tau + \left( \frac{\partial U}{\partial V} \right)_\tau dV + pdV \equiv C_V d\tau + \left( \frac{\partial U}{\partial V} \right)_\tau dV + pdV \end{aligned}$$

with  $U = U(\tau, V)$ .

For (calorically) perfect gas,  $U = U(\tau)$  and so

$$\begin{aligned} d\sigma &= C_V \frac{d\tau}{\tau} + \frac{N}{V} dV = C_V \frac{d\tau}{\tau} + \frac{-N d\rho}{\rho} \\ \implies \sigma_2 - \sigma_1 &= C_V \ln \left( \frac{\tau_2}{\tau_1} \right) + N \ln \left( \frac{V_2}{V_1} \right) = \ln \left[ \left( \frac{\tau_2}{\tau_1} \right)^{C_V} \left( \frac{V_2}{V_1} \right)^N \right] \\ \implies \exp(\sigma_2 - \sigma_1) &= \left[ \left( \frac{\tau_2}{\tau_1} \right)^{\frac{C_V}{N}+1} \frac{p_1}{p_2} \right]^N \quad \text{or} \quad \frac{p_2}{p_1} = \left( \frac{\tau_2}{\tau_1} \right)^{\frac{\gamma}{\gamma-1}} \exp \left[ \frac{\sigma_2 - \sigma_1}{-N} \right] \end{aligned}$$

Starting from Sec. 2.16 *The Channel Flow Equations* of Ch. 2 of Oates (1997) [1], pp. 39-40, consider how heating  $Q'$  inside an (control volume) element of fluid in a duct changes the (specific) enthalpy  $h'$ :

$$Q' \in \Omega^1(\Sigma) \otimes \Omega^n(N)$$

$$h' \in \Omega^1(\Sigma) \otimes \Omega^{n-1}(N)$$

where it is now clear that the integration over the (control volume) element of the fluid on spatial manifold  $N$  is distinguished from the integration (usually along a curve representing a thermodynamic process) over the manifold  $\Sigma$  of *thermodynamic states*.

By energy conservation,

$$Q' = \mathbf{d}h'$$

where  $\mathbf{d}$  is the exterior derivative *on*  $N$ , not on  $\Sigma$ .  $h'$ , the “specific enthalpy”, is a *Bernoulli invariant* of the flow<sup>2</sup>.

$h'$  is also called *stagnation enthalpy*.

For instance, integrate over the control volume, and then over a 1-dimensional flow (cylindrical volume element)

$$\int_{V_0} Q' = \int_{V_0} \mathbf{d}h' = (h')_2 A_2 - (h')_1 A_1 = (h_0 + \frac{u^2}{2})_2 \rho_2 A_2 - (h_0 + \frac{u^2}{2})_1 \rho_1 A_1$$

But one also have to account for the thermodynamic process during the flow across the element control volume, on  $\Sigma$ :

$$\tau d\sigma = dU + pdV = dH - V dp = dH - \frac{dp}{\rho/NM}$$

So for the Bernoulli invariant  $h'$ ,  $h'$  has a piece  $h \equiv H/NM \in C^\infty(\Sigma)$ , that lives on  $\Sigma$ .

In conclusion,

$$(18) \quad Q' = \mathbf{d} \left( \frac{u^2}{2} \right) + \mathbf{d}h = \mathbf{d} \left( \frac{u^2}{2} \right) + \frac{\tau}{NM} d\sigma + \frac{dp}{\rho} \equiv \mathbf{d} \left( \frac{u^2}{2} \right) + \tau ds + \frac{dp}{\rho}$$

<sup>2</sup>“Euler equations (fluid dynamics)” *Wikipedia*, [https://en.wikipedia.org/wiki/Euler\\_equations\\_\(fluid\\_dynamics\)](https://en.wikipedia.org/wiki/Euler_equations_(fluid_dynamics))

which is Eq. (2.61) of Oates (1997) [1].

Consider mass conservation for fluid flow across a control volume element:

$$\dot{m} = \frac{d}{dt} \int_{V_0} \rho \text{vol}^n = \int_{V_0} \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{V_0} di_u \rho \text{vol}^n = 0 + \int_{\partial V_0} i_u \rho \text{vol}^n = (\rho u A)_2 - (\rho u A)_1$$

Then for  $\rho u A = \text{const.}$ ,  $\frac{d\rho}{\rho} + \frac{du}{u} + \frac{dA}{A} = 0$ .

Assuming an ideal gas,  $p = \frac{N\tau}{V} = \frac{\rho\tau}{M}$

$$\frac{dA}{A} + \frac{du}{u} + \frac{d\rho}{\rho} = \frac{dA}{A} + \frac{du}{u} + \frac{d(pM/\tau)}{\rho}$$

Then for  $\frac{d(pM/\tau)}{\rho}$  only, and using Eq. 18, which is from energy conservation where  $Q'$  heating (or heat dissipation) is also included,

$$\begin{aligned} \frac{d(pM/\tau)}{\rho} &= \frac{M}{\tau} \frac{dp}{\rho} + \frac{pM}{-\tau^2} \frac{d\tau}{\rho} = \\ &= \left( Q' - \mathbf{d} \left( \frac{u^2}{2} \right) - \frac{\tau d\sigma}{NM} \right) \frac{M}{\tau} + \frac{pM}{-\tau^2} \frac{d\tau}{\rho} \end{aligned}$$

The last term,  $\frac{pM}{-\tau^2} \frac{d\tau}{\rho}$ , is resolved by considering again the Bernoulli invariant  $h$ , and the fact that it remains constant along the flow:

$$\begin{aligned} \frac{pM}{-\tau^2} \frac{d\tau}{\rho} &= \frac{pMV}{-\tau^2 MN} d\tau = -\frac{d\tau}{\tau} \quad (\text{ideal gas}) \\ h &= h_1 + \frac{u_1^2}{2} = h_2 + \frac{u_2^2}{2} \implies dh = C_p/(MN) d\tau = u du \\ \implies \frac{d\tau}{\tau} &= \frac{u du}{(C_p/MN)\tau} = \frac{u du}{\left( \frac{\gamma}{M(\gamma-1)} \right) \tau} = M \frac{\gamma-1}{\gamma} \frac{u du}{\tau} \end{aligned}$$

Thus

$$Q' \frac{M}{\tau} - \frac{d\sigma}{N} - \frac{M}{\gamma} \frac{u du}{\tau} = Q' \frac{M}{\tau} - \frac{d\sigma}{N} - \mathfrak{M}^2 \frac{du}{u}$$

with Mach number  $\mathfrak{M} := \frac{u}{a}$  and speed of sound  $a^2 = \frac{\gamma\tau}{M}$  which comes from the adiabatic process of compressing and the expansion of gas longitudinally  $\left( \frac{\partial p}{\partial \rho} \right)_{\text{adiabatic}} = a^2$ .

And so for fluid (ideal gas) flow through a (control volume, differential) element of a (area) duct, including heating (or heat dissipation),

$$(19) \quad \frac{dA}{A} + (1 - \mathfrak{M}^2) \frac{du}{u} = \frac{d\sigma}{N} - Q' \frac{M}{\tau}$$

6.0.1. *Adiabatic Flow of Ideal Gas.* Let  $Q' = 0$  for adiabatic flow (by definition of “adiabatic”). Then from Eq. [19](#)

$$\frac{dA}{A} + (1 - \mathfrak{M}^2) \frac{du}{u} = \frac{d\sigma}{N}$$

Now  $d\sigma \geq 0$ , always (2nd. law of Thermodynamics).  $N > 0$ .

$dA < 0$  for *converging* nozzle.  $dA > 0$  for *diverging* nozzle.

For isentropic flow  $d\sigma = 0$ ,

for  $dA < 0$ , for  $\mathfrak{M} < 1$ ,  $du > 0$  (flow is accelerating)

for  $\mathfrak{M} > 1$ ,  $du < 0$  (flow is decelerating)

At the throat,  $dA = 0$ ,  $\mathfrak{M} = 1$ .

However, for  $d\sigma > 0$ , at the throat,  $dA = 0$ , for acceleration  $du > 0$ ,  $\mathfrak{M} < 1$  and for deceleration  $du < 0$ ,  $\mathfrak{M} > 1$ . So from subsonic flow in the initial part (usually the converging part), then the “effective throat” of the nozzle is shifted slightly downstream.

cf. Ch. 2 Problems, Oates (1997) [\[1\]](#)

**Problem 2.1.**

Starting from the thermodynamic identity,

$$d\sigma = \frac{dU}{\tau} + \frac{p}{\tau} dV = \frac{dU}{\tau} + \frac{NdV}{V}$$

$U = U(\tau)$  for an (perfect) ideal gas, and  $C_V = \left(\frac{\partial U}{\partial \tau}\right)_V = \frac{dU}{d\tau}$  so

$$d\sigma = C_V \frac{d\tau}{\tau} + \frac{N}{V} dV$$

For an isentropic process,  $\int d\sigma = 0$ , so

$$\int d\sigma = 0 = \int C_V \frac{d\tau}{\tau} + N \ln \frac{V_2}{V_1}$$

For  $C_V = A + B\tau + C\tau^2 + D\tau^3$ ,

$$\begin{aligned} N \ln \frac{V_1}{V_2} &= N \ln \frac{\rho_2}{\rho_1} = A \ln \left( \frac{\tau_2}{\tau_1} \right) + B(\tau_2 - \tau_1) + \frac{C(\tau_2 - \tau_1)^2}{2} + D \frac{(\tau_2 - \tau_1)^3}{3} \\ \implies \frac{\rho_2}{\rho_1} &= \left( \left( \frac{\tau_2}{\tau_1} \right) \exp \left[ \frac{B}{A}(\tau_2 - \tau_1) + \frac{C}{2A}(\tau_2 - \tau_1)^2 + \frac{D}{3A}(\tau_2 - \tau_1)^3 \right] \right)^{A/N} \end{aligned}$$

Over a “wide” temperature range,  $C_p$  is usually given as a polynomial, from empirically measured parameters.

For  $C_p$ , consider

$$dH = dU + pdV + Vdp = \tau d\sigma + Vdp = \left( \frac{\partial H}{\partial \tau} \right)_p d\tau + \left( \frac{\partial H}{\partial p} \right)_\tau dp$$

$$Q = dU + pdV = dH - Vdp$$

$$C_p = \left( \frac{\partial H}{\partial \tau} \right)_p$$

$$Q = \tau d\sigma = dH - Vdp \text{ or } d\sigma = \frac{dH}{\tau} - \frac{V}{\tau} dp = \frac{dH}{\tau} - \frac{Ndp}{p}$$

$$\implies \int d\sigma = \int \frac{C_p d\tau}{\tau} - N \ln \frac{p_2}{p_1}$$

$$\implies \frac{p_2}{p_1} = \left( \left( \frac{\tau_2}{\tau_1} \right) \exp \left[ \frac{B}{A}(\tau_2 - \tau_1) + \frac{C}{2A}(\tau_2 - \tau_1)^2 + \frac{D}{3A}(\tau_2 - \tau_1)^3 \right] \right)^{A/N}$$

cf. Ch. 3 Chemical Rockets, Oates (1997) [\[1\]](#)

rocket volume  $R_0$

$\partial R_0$  includes

internal surfaces wetted by fluid (chamber, pipes, pumps, etc.)

Consider rocket on thrust stand. Consider force on thrust stand.

$\sum_0$  outer surface area.

$\sum_c$  inner chamber surface area.

Now consider the total mass of the gas propellant, in the chamber and nozzle, represented by a differential  $n$ -form on (spatial) manifold  $N$ , representing spatial points, of dimension  $\dim N = n$ , denoted by  $m$ :

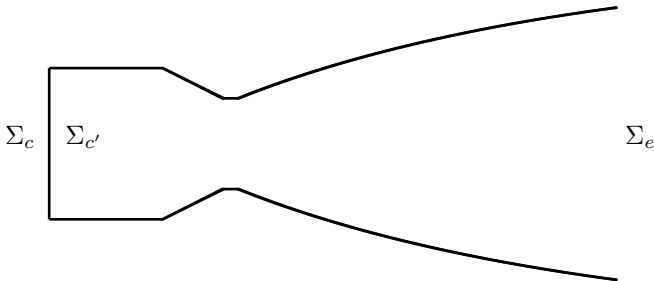
$$m = \int_{B(t)} \rho \text{vol}^n$$

where  $B(t)$  is the control volume, as a function of time  $t$ , and  $\text{vol}^n \in \Omega^n(N)$  is the volume form,  $\Omega^n(N)$  is the set of all  $n$ -forms on  $N$ .

Then, by differential geometry,

$$\dot{m} = \int_{B(t)} \left[ \frac{\partial \rho}{\partial t} \text{vol}^n + \mathcal{L}_u \rho \text{vol}^n \right] = \int_{B(t)} \left[ \frac{\partial \rho}{\partial t} \text{vol}^n + di_u \rho \text{vol}^n \right] = \int_{B(t)} \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial B} i_u \rho \text{vol}^n = \int_{B(t)} \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial B} \rho u^i dS_i$$

where the last equality was obtained by Stoke’s law.



Consider the total momentum, the sum of all the momentums of the fluid particles,  $\Pi$ , and its time derivative,  $\dot{\Pi}$ ,

$$\Pi = \int_{B(t)} \rho u^i \text{vol}^n \otimes e_i$$

$$\dot{\Pi} = \int_{B(t)} \frac{\partial(\rho u^i)}{\partial t} \text{vol}^n \otimes e_i + \int_{B(t)} d(\rho u^i i_u \text{vol}^n) \otimes e_i = \int_{B(t)} \frac{\partial(\rho u^i)}{\partial t} \text{vol}^n \otimes e_i + \int_{\partial B(t)} \rho u^i i_u \text{vol}^n \otimes e_i$$

For steady state  $\frac{\partial(\rho u^i)}{\partial t} = 0$ ,  $\dot{\Pi} = \int_{\partial B(t)} \rho u^i i_u \text{vol}^n \otimes e_i$ .

Now

$$\partial B(t) = A_e \coprod \Sigma_{c'}$$

$\Sigma_{c'} \equiv$  internal surface of chamber (orientation, or outward normal, chosen to point into rocket stages, i.e. outward from engine).

Suppose  $u = 0$  on  $\Sigma_{c'}$ . Then

$$\dot{\Pi} = \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i$$

Consider the external forces on a fluid. The external forces on a fluid, acting on the surface of the fluid, is “wrapped up” in the stress tensor  $T = T^{ij} e_i \otimes e_j \in TN \otimes TN$ . On a surface,

$$\int T^{ij} dS_j \otimes e_i$$

and so the external forces on a fluid, decomposing  $T$ , is

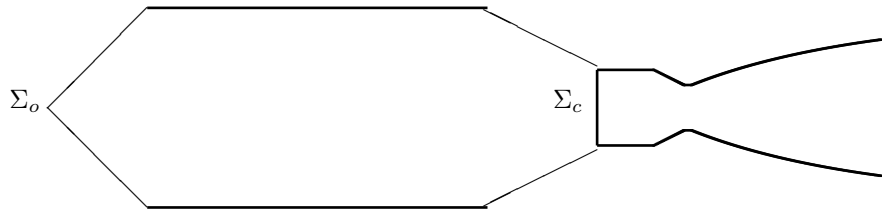
$$\begin{aligned}
 (20) \quad & \int_{\Sigma_{c'} + A_e} T^{ij} dS_j \otimes e_i = \\
 & = - \int_{\Sigma_{c'} + A_e} (p - p_a) g^{ij} dS_j \otimes e_i + (\text{visc})_{\Sigma_{c'}} + (\text{visc})_{A_e} = \\
 & = \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i
 \end{aligned}$$

where in the last equality, we equate the sum of the external forces on the fluid, to the kinematics of the external fluid, which is what Newton's second law tells us to do.

Consider the  $-\int_{\Sigma_{c'} + A_e} (p - p_a) g^{ij} dS_j \otimes e_i$  term. Consider a gas bubble in water. For a gas bubble of high pressure  $p$ , expanding, after, say, an explosion underwater, then  $-p g^{ij} dS_j \otimes e_i$  is from Newton's 3rd. law, the pressure pushing back on expanding gas bubble, pushing inward ( $-dS_j$ ).  $p_a g^{ij} dS_j \otimes e_i$  is due to the ambient pressure, pushing outward. In the context of the gas propellant in the rocket nozzle, immersed in the atmosphere, (the atmospheric) air pushes outward to try to disperse the gas propellant. If  $p_a = 0$  (in vacuum), there's no other gas molecules around; the gas bubble drifts if  $p = 0$ , or expands out unrestrained, if  $p \neq 0$  or  $p > 0$  (defined how I defined it,  $p \geq 0$  always).

Rearranging Eq. 20,

$$(21) \quad \int_{\Sigma_{c'}} T^{ij} dS_j \otimes e_i = - \int_{A_e} T^{ij} dS_j \otimes e_i + \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i = \int_{A_e} (p - p_a) g^{ij} dS_j \otimes e_i - (\text{visc})_{\Sigma_{c'}} - (\text{visc})_{A_e} + \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i$$



Now consider the (sum of all external) forces on the rocket (i.e. the vehicle, rocket, and the “interesting part”, payload), and plugging in Eq. 21, but not including body forces,

$$\begin{aligned}
 (22) \quad \sum_{\text{rocket}} F_{\text{ext}} &= \int_{\Sigma_o + \Sigma_c} T^{ij} dS_j \otimes e_i = \int_{\Sigma_o} T^{ij} dS_j \otimes e_i + \int_{\Sigma_c} T^{ij} dS_j \otimes e_i = \int_{\Sigma_o} T^{ij} dS_j \otimes e_i - \int_{\Sigma_{c'}} T^{ij} dS_j \otimes e_i = \\
 &= - \int_{\Sigma_o} p_a g^{ij} dS_j \otimes e_i + (\text{visc})_{\Sigma_o} - \left[ \int_{A_e} (p - p_a) g^{ij} dS_j \otimes e_i - (\text{visc})_{\Sigma_{c'}} - (\text{visc})_{A_e} + \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i \right] = \\
 &= - \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i - \int_{A_e} (p - p_a) g^{ij} dS_j \otimes e_i + (\text{visc})_{\Sigma_o} - \int_{\Sigma_o} p_a g^{ij} dS_j \otimes e_i + (\text{visc})_{\Sigma_{c'}} + (\text{visc})_{A_e}
 \end{aligned}$$

Compare this to Eq. (3.1) or Eq. (3.3) on pp. 64, Ch. 3 of Oates (1997) [1]:

$$\mathbf{F} = - \iint_{A_e} (\rho \mathbf{u}) \mathbf{u} \cdot d\mathbf{s} - \iint_{A_e} (p - p_a) d\mathbf{s} + (\text{visc})_{\Sigma_o} - \iint_{\Sigma_o} (p - p_a) d\mathbf{s}$$

I agree with the first 3 terms of Eq. (3.3) on pp. 64 of Oates (1997), with Eq. 22. I don't agree with Oates on the fourth term. The fourth term represents form drag, which depends on the shape of the external, outward rocket shell, and the atmosphere.

Oates says that it “represents the effect of the pressure imbalance on the external surface and is termed the form drag.” However, it is unclear to me what the  $p$  is for Oates [1] other than the ambient atmospheric pressure. And, if we're in vacuum, form drag should go to 0. So for  $p_a = 0$ , then form drag is zero.

In summary,

$$\begin{aligned}
 (\text{visc})_{\Sigma_o} &=: \text{skin drag} \\
 - \int_{\Sigma_o} p_a g^{ij} dS_j \otimes e_i &=: \text{form drag}
 \end{aligned}$$

Term  $(\text{visc})_{A_e}$  of Eq. ?? is neglected, which I agree with.  $(\text{visc})_{\Sigma_{c'}}$  was subtracted away when reversing orientation. I do not agree with this step. What remains invariant is stress tensor  $T$ , which should be orientation preserving. The shear force, in off-diagonal entries of  $T$ , include viscosity. Because clearly, in vector form,

$$\int (\mathbf{T} \cdot \mathbf{n}) dA = \int (T(\cdot, \mathbf{n}) dA \xrightarrow{\mathbf{n} \rightarrow -\mathbf{n}} \int (T \cdot (-\mathbf{n})) dA = - \int T \cdot \mathbf{n} dA$$

Nevertheless, the thrust, denoted  $F_T$ , is

$$(23) \quad F_T = M(t) \mathbf{a} = - \int_{A_e} \rho u^i i_u \text{vol}^n \otimes e_i - \int_{A_e} (p - p_a) g^{ij} dS_j \otimes e_i$$

Assume  $u^i = u_{\text{exh}}$  on  $A_e$  (constant), and  $p, p_a$  uniform (constant) on  $A_e$ ,

$$(24) \quad F_T \equiv T = \dot{m} u_e + (p_{\text{exh}} - p_a) A_e$$

Compare this to Eq. (2-14) of Biblarz and Sutton (2001) [3].

Let  $T = \dot{m} u_c$ , where  $u_c$  is the effective velocity (cf. Eq. 6).

$$T = \dot{m} u_c = \dot{m} u_{\text{exh}} + (p_{\text{exh}} - p_a) A_e$$

so that

$$(25) \quad u_c = u_{\text{exh}} + (p_{\text{exh}} - p_a) A_e / \dot{m}$$

Compare this to Eq. (2-16) of Biblarz and Sutton (2001) [3].

Biblarz and Sutton (2001) [3] then considers the following special case and approximation: approximate  $u_c = u_{\text{exh}}$  (due to the special case that  $p_{\text{exh}} = p_a$ , approximately), then

$$(26) \quad F_T = \dot{m} u_c$$

Compare this to Eq. (2-17) of Biblarz and Sutton (2001) [3].

6.0.2. *Effective exhaust velocity, Specific impulse  $I_{sp}$ .* From Eq. 25,

$$\begin{aligned}
 (27) \quad u_c &= u_{\text{exh}} + (p_{\text{exh}} - p_a) A_e / \dot{m} = u_{\text{exh}} \left( 1 + \frac{(p_{\text{exh}} - p_a) A_e}{u_{\text{exh}} \dot{m}} \right) = u_{\text{exh}} \left( 1 + \frac{(p_{\text{exh}} - p_a)}{u_{\text{exh}} (\rho_{\text{exh}} u_{\text{exh}})} \right) = \\
 &= u_{\text{exh}} \left( 1 + \frac{p_{\text{exh}}}{\rho_{\text{exh}} u_{\text{exh}}^2} \left( 1 - \frac{p_a}{p_{\text{exh}}} \right) \right) = u_{\text{exh}} \left( 1 + \frac{1}{\gamma \mathfrak{M}_{\text{exh}}^2} \left( 1 - \frac{p_a}{p_{\text{exh}}} \right) \right)
 \end{aligned}$$

$\mathfrak{M} := \frac{u}{a}$ , speed of sound  $a^2 = \frac{\gamma \tau}{M}$ ,  $pV = N\tau$  or  $p = \rho \frac{\tau}{M}$ ;  $\mathfrak{M}_{\text{exh}}^2 = \frac{u_{\text{exh}}^2}{a^2} = \frac{u_{\text{exh}}^2 \rho_{\text{exh}}}{\gamma p_{\text{exh}}}$ .

On Section 3.3 Acceleration of a Rocket, pp. 66, of Oates (1997) [1], the acceleration of a rocket, in general, is

$$M \mathbf{a} = - \int_{A_e} u^i \rho i_u \text{vol}^n \otimes e_i - \int_{A_e} (p - p_a) g^{ij} dS_j \otimes e_i - Mg - D$$

In direction of flight,

$$\begin{aligned}
 M \frac{dv}{dt} &= \dot{m} u_{\text{exh}} + (p_{\text{exh}} - p_a) A_e - Mg \cos \theta - D \text{ or} \\
 \frac{dv}{dt} &= \frac{\dot{m} u_{\text{exh}}}{M} + \frac{(p_{\text{exh}} - p_a) A_e}{M} - g \cos \theta - \frac{D}{M}
 \end{aligned}$$

Let the mass of the rocket be

$$\begin{aligned} M &= M(t) = M(0) - m(t) \text{ s.t.} \\ m(0) &= 0 \\ m(t_p) &= m_p \text{ so} \\ M(0) &= M_0 + m_p \end{aligned}$$

Now let

$$\begin{aligned} \lambda &= \frac{m_L}{m_0} & \text{payload ratio} &\equiv \lambda & \lambda &= \frac{m_L}{M(0)} \\ \delta &= \frac{m_d}{m_0} & \text{dead weight ratio} &\equiv \delta & \delta &= \frac{m_d}{M(0)} \end{aligned}$$

and so, in terms of the payload ratio  $\lambda$  and dead weight ratio  $\delta$ ,

$$\frac{M(0)}{M(0) - m_p} = \frac{M(0)}{m_d + m_L} = \frac{1}{\delta + \lambda}$$

Assume constant burn rate:  $M(t) = M(0) - \frac{m_p}{t_f}t$ .

Free space case:  $g = D = 0$

$$\begin{aligned} \frac{dv}{dt} &= \frac{-\dot{M}u_{\text{exh}}}{M} + \frac{(p_{\text{exh}} - p_a)A_e}{M} \\ \Delta v_{\text{tot}} &= -u_{\text{exh}} \ln \left( \frac{M(0) - m_p}{M(0)} \right) + \frac{(p_{\text{exh}} - p_a)A_e}{\frac{-m_p}{t_f}} \ln \left( M(0) - \frac{m_p}{t_f}t \right) \Bigg|_0^{t_f} = \\ &= -u_{\text{exh}} \ln \left( \frac{M(0) - m_p}{M(0)} \right) + \frac{(p_{\text{exh}} - p_a)A_e}{\frac{-m_p}{t_f}} \ln \left( \frac{M(0) - m_p}{M(0)} \right) = \\ &= \left[ u_{\text{exh}} + \frac{(p_{\text{exh}} - p_a)A_e}{\dot{m}} \right] \ln \left( \frac{M(0)}{M(0) - m_p} \right) \end{aligned}$$

For an atmosphere-free planet,  $D = 0$ ,

$$\begin{aligned} \dot{v} &= \frac{-\dot{M}u_{\text{exh}}}{M} + \frac{(p_{\text{exh}} - p_a)A_e}{M} - g \cos \theta \\ \Rightarrow \Delta v &= \left( u_{\text{exh}} + \frac{(p_{\text{exh}} - p_a)A_e}{\dot{m}} \right) \ln \left( \frac{M(0)}{M(0) - m_p} \right) - g \cos \theta t_f \end{aligned}$$

## 7. THERMOCHEMISTRY, THERMODYNAMICS, ONE-DIMENSIONAL FLUID FLOW

### 8. MULTIPLE STAGING

cf. pp. 67, **Multiple-Stage Rockets**, Ch. 3, Oates (1997) [1]

### 9. LAGRANGIAN

9.1. **Review.** Consider a solid cylinder with surface area  $A_{\text{cyl}} = A_{\text{sides}} \amalg A_{\text{top}} \amalg A_{\text{bottom}}$  s.t.

$$\begin{aligned} dA_{\text{sides}} &= ad\varphi dz & \hat{n} &= \frac{x\mathbf{e}_x + y\mathbf{e}_y}{\|x\mathbf{e}_x + y\mathbf{e}_y\|} = \mathbf{e}_r = \cos \varphi \mathbf{e}_x + \sin \varphi \mathbf{e}_y \\ dA_{\text{top}} &= dA_{\text{bottom}} = 2\pi r dr & \hat{n} &= \pm \mathbf{e}_z \quad (\text{for top and bottom, respectively}) \end{aligned}$$

for a fixed, *body-frame* coordinates  $\mathbf{e}_x, \mathbf{e}_y, \mathbf{e}_z$ .

Let

$$A_{\text{projected}} = \int_A (\cos \beta) dA \chi_{\cos \beta \geq 0}$$

where  $A$  is the original area

$\beta$  is the angle between normal to surface  $A$  and normal to arbitrary plane onto which we project and

$\chi_{\cos \beta \geq 0}$  is an indicator function, a fancy way/terminology for saying s.t.  $\chi_{\cos \beta \geq 0} = \begin{cases} 1 & \text{if } \cos \beta \geq 0 \\ 0 & \text{otherwise} \end{cases}$ .

Consider  $\frac{\mathbf{u}}{\|\mathbf{u}\|} \equiv \hat{u} = \hat{u}_x \mathbf{e}_x + \hat{u}_y \mathbf{e}_y + \hat{u}_z \mathbf{e}_z$ .

Consider now these sanity checks: consider  $\hat{u} = \mathbf{e}_z$ .

$$\begin{aligned} \int_{A_{\text{top}}} \hat{u} \cdot \hat{n} 2\pi r dr &= 2\pi \int_0^a r dr = \pi a^2 \\ \int_{A_{\text{bottom}}} \hat{u} \cdot \hat{n} 2\pi r dr \chi_{\hat{u} \cdot \hat{n} \geq 0} &= 0 \text{ since } \hat{u} \cdot \hat{n} = -1 \\ \int_{A_{\text{side}}} \hat{u} \cdot \hat{n} 2\pi r dr &= 0 \text{ since } \hat{u} \cdot \hat{n} = 0 \\ \Rightarrow A_{\text{projected}} &= \pi a^2 \end{aligned}$$

Consider  $\hat{u} = \mathbf{e}_x$  or  $= \mathbf{e}_y$  or  $= \mathbf{e}_r$ .  $\hat{u} \cdot \hat{n} = 0$  on  $A_{\text{top}}$  and on  $A_{\text{bottom}}$ .

$$\int_{A_{\text{side}}} \hat{u} \cdot \hat{n} a d\varphi dz \chi_{\hat{u} \cdot \hat{n} \geq 0} = ah \int_{-\pi}^{\pi} d\varphi \cos \varphi \chi_{\cos \varphi \geq 0} = ah \int_{-\pi/2}^{\pi/2} d\varphi \cos \varphi = 2ah$$

Indeed, this result is rotationally symmetric about  $\mathbf{e}_z$ , for if  $\hat{u} = \mathbf{e}_y$ .

$$\int_{A_{\text{side}}} \sin \varphi d\varphi \chi_{\sin \varphi \geq 0} = \int_0^{\pi} d\varphi \sin \varphi = -\cos \varphi \Big|_0^{\pi} = 2$$

and so  $A_{\text{projected}} = 2ah$ .

Consider  $\hat{u} = \hat{u}_\rho \mathbf{e}_\rho + \hat{u}_z \mathbf{e}_z$ . Then  $A_{\text{projected}} = \hat{u}_z \pi a^2 + \hat{u}_\rho 2ah$ .

Returning back to the drag

$$\mathbf{F}_d = \frac{-1}{2} \rho C_a A |\mathbf{u}| \mathbf{u}$$

The components of the drag force in 2-dim. Cartesian coordinates are

$$(\mathbf{F}_d)_1 = \frac{-1}{2} \rho C_a A (\dot{x}_1^2 + \dot{x}_2^2)^{1/2} \dot{x}_1$$

$$(\mathbf{F}_d)_2 = \frac{-1}{2} \rho C_a A (\dot{x}_1^2 + \dot{x}_2^2)^{1/2} \dot{x}_2$$

In spherical coordinates,

$$Q_r = F_1 \frac{\partial x_1}{\partial r} + F_2 \frac{\partial x_2}{\partial r} = \frac{-1}{2} \rho C_a A (\dot{r}^2 + (r\dot{\varphi})^2)^{1/2} \dot{r}$$

$$Q_\varphi = F_1 \frac{\partial x_1}{\partial \varphi} + F_2 \frac{\partial x_2}{\partial \varphi} = \frac{-1}{2} \rho C_a A (\dot{r}^2 + (r\dot{\varphi})^2)^{1/2} r^2 \dot{\varphi}$$

since  $(\dot{r}c_\varphi - rs_\varphi\dot{\varphi})c_\varphi + (\dot{r}s_\varphi + rc_\varphi\dot{\varphi})s_\varphi = \dot{r}$  (cf. Chapter 6 “Lagrangian Dynamics” <sup>3</sup>)

$$(\dot{r}c_\varphi - rs_\varphi\dot{\varphi})(-s_\varphi) + (\dot{r}s_\varphi + rc_\varphi\dot{\varphi})c_\varphi = r\dot{\varphi}$$

Clearly, the so-called “Rayleigh’s dissipation function” (Eq. (1.67) on pp. 23 of Goldstein, et. al. [11]) is

$$\mathcal{F} = \frac{1}{3} \rho C_a A (\dot{x}_1^2 + \dot{x}_2^2)^{3/2}$$

One can also write down a generalized force for the thrust

$$F_{\text{thrust}} = T = \dot{m}_p \mathbf{u}_e + (p_e - p_a) A_e \mathbf{e}_z$$

<sup>3</sup><http://ice.as.arizona.edu/~dpsaltis/Phys422/chapter6.pdf>

where  $\mathbf{u}_e$  is in the direction opposite of the  $z$  axis of symmetry of the *body-frame*, as with  $\mathbf{e}_z$ .

One should also consider  $\mathbf{e}_z$  to be a vector in the *body-frame* that's arbitrary to represent cold-gas thrusters:

$$F_{\text{thrust}} \equiv T = (\dot{m}_p u_e + (p_e - p_a) A_e) \mathbf{e}_z$$

## 10. MISSION DESIGN

### 10.1. Motion in a central field.

10.1.1. *Reduced Mass.* cf. Sec. 30, *The reduced mass*, Landau and Lifshitz (1976) [5])

Lagrangian of system of 2 interacting particles, with potential energy of interaction of 2 particles depending only on distance between them, i.e. magnitude of difference in their radius vectors:

$$(28) \quad \mathcal{L} = \frac{1}{2} m_1 \dot{\mathbf{r}}_1^2 + \frac{1}{2} m_2 \dot{\mathbf{r}}_2^2 - U(|\mathbf{r}_1 - \mathbf{r}_2|)$$

Let  $\mathbf{r} \equiv \mathbf{r}_1 - \mathbf{r}_2$  be the relative position vector.

Let origin be at the center of mass  $\mathbf{r}_{\text{cm}}$ . In general,

$$(29) \quad \mathbf{r}_{\text{cm}} = \frac{m_1 \mathbf{r}_1 + m_2 \mathbf{r}_2}{m_1 + m_2}$$

Then

$$\mathbf{r}_{\text{cm}} = \frac{m_1 \mathbf{r}_1 + m_2 (\mathbf{r}_2 - \mathbf{r})}{m_1 + m_2} \implies \begin{aligned} \mathbf{r}_1 &= \mathbf{r}_{\text{cm}} + \frac{m_2 \mathbf{r}}{m_1 + m_2} \\ \mathbf{r}_2 &= \mathbf{r}_{\text{cm}} + \frac{m_1 \mathbf{r}}{m_1 + m_2} \end{aligned}$$

With

$$\frac{1}{2} m_1 \dot{\mathbf{r}}_1^2 = \frac{1}{2} m_1 \left( \dot{\mathbf{r}}_{\text{cm}}^2 + \frac{m_2}{M} \dot{\mathbf{r}}_{\text{cm}} \cdot \dot{\mathbf{r}} + \frac{m_2^2}{M^2} \dot{\mathbf{r}}^2 \right)$$

Then

$$(30) \quad \mathcal{L} = \frac{1}{2} M \dot{\mathbf{r}}_{\text{cm}}^2 + m \dot{\mathbf{r}}_{\text{cm}} \cdot \dot{\mathbf{r}} + \frac{1}{2} \frac{m}{M} \dot{\mathbf{r}}^2 - U(r)$$

where *reduced mass* is given by

$$(31) \quad \text{reduced mass } \mu \equiv m = \frac{m_1 m_2}{m_1 + m_2} = \frac{m_1 m_2}{M}$$

Thus the problem of motion of 2 interacting particles is equivalent to that of motion of 1 particle in a given external field  $U(r)$ .

Problem: system consists of 1 particle of mass  $M$ , and  $n$  particles with equal masses  $m$ . (cf. pp. 30, Landau and Lifshitz (1976) [5]).

Let  $\mathbf{R} \equiv$  radius vector particle of mass  $M$ ,

$\mathbf{R}_\alpha$ ,  $\alpha = 1, 2, \dots, n$  radius vectors of particles of mass  $m$ ,

$\mathbf{r}_\alpha \equiv \mathbf{R}_\alpha - \mathbf{R}$  (just define this vector), then

$$\mathbf{r}_{\text{cm}} = \frac{M \mathbf{R} + \sum_{\alpha=1}^n m_\alpha \mathbf{R}_\alpha}{M + \sum_{\alpha=1}^n m_\alpha} = \frac{M \mathbf{R} + m \sum_{\alpha=1}^n (\mathbf{r}_\alpha + \mathbf{R})}{M_{\text{tot}}} = \mathbf{R} + \frac{m}{M_{\text{tot}}} \sum_{\alpha=1}^n \mathbf{r}_\alpha$$

$$\mathbf{R} = \mathbf{r}_{\text{cm}} - \frac{m}{M_{\text{tot}}} \sum_{\alpha=1}^n \mathbf{r}_\alpha$$

$$\mathbf{R}_\alpha = \mathbf{r}_\alpha + \mathbf{R}$$

and so for the kinetic energy part of

$$\mathcal{L} = \frac{1}{2} M \dot{\mathbf{R}}^2 + \frac{1}{2} m \sum \dot{\mathbf{R}}_\alpha^2 - U$$

consider

$$M V^2 + m \sum (\mathbf{v}_\alpha + \mathbf{V})^2 = M V^2 + m \sum (v_\alpha^2 + 2 \mathbf{v}_\alpha \cdot \mathbf{V} + V^2) = M_{\text{tot}} V^2 + m \sum v_\alpha^2 + 2m \sum \mathbf{v}_\alpha \cdot \mathbf{V} =$$

$$M_{\text{tot}} \left( v_{\text{cm}}^2 - 2 \frac{m}{M_{\text{tot}}} \mathbf{v}_{\text{cm}} \cdot \sum_{\alpha=1}^n \mathbf{v}_\alpha + \frac{m^2}{M_{\text{tot}}^2} \left( \sum_{\alpha=1}^n v_\alpha \right)^2 \right) + m \sum v_\alpha^2 + 2m \sum \mathbf{v}_\alpha \cdot \left( \mathbf{v}_{\text{cm}} - \frac{m}{M_{\text{tot}}} \sum \mathbf{v}_\alpha \right) =$$

$$M_{\text{tot}} v_{\text{cm}}^2 - \frac{m^2}{M_{\text{tot}}} \left( \sum_{\alpha=1}^n v_\alpha \right)^2 + m \sum v_\alpha^2$$

Thus

$$\mathcal{L} = \frac{1}{2} \left( M_{\text{tot}} v_{\text{cm}}^2 - \frac{m^2}{M_{\text{tot}}} \left( \sum_{\alpha=1}^n v_\alpha \right)^2 + m \sum v_\alpha^2 \right) - U$$

### Part 4. AE121

Let's translate  $\frac{\gamma}{\gamma-1}$  between physicists and engineers:

$$(32) \quad \boxed{\frac{\gamma}{\gamma-1} = \frac{\frac{C_p}{C_V}}{\frac{N}{C_V}} = \frac{C_p}{N} = \frac{c_p M N}{N} = \frac{c_p k_B}{R}}$$

10.1.2. *Speed of sound.* From pp. 179, Problem 10 “Isnetropic relations of ideal gas” of Chapter 6: Ideal Gas of Kittel and Kroemer [8], recall isentropic bulk moduli  $B_\sigma$

$$B_\sigma := -v \left( \frac{\partial p}{\partial V} \right)_\sigma = \gamma \frac{p_i V_i^\gamma}{V^\gamma} = \gamma p$$

with  $p = \frac{p_i V_i^\gamma}{V^\gamma}$ .

Very little heat transfer in sound wave. For velocity (magnitude) i.e. speed of sound,  $a$

$$a = \left( \frac{B\sigma}{\rho} \right)^{1/2} = \left( \frac{\gamma p}{\rho} \right)^{1/2}$$

Now

$$p = \frac{N\tau}{V}$$

$$\frac{p}{\rho} = \frac{N\tau}{V} \frac{1}{\left(\frac{MN}{V}\right)} = \frac{\tau}{M}$$

Now  $p = \rho R T$  (outside of theoretical physics, people use the so-called universal gas constant  $R$ ).

$$p = \rho R T = \frac{N\tau}{V} = \frac{MN}{V} R \frac{\tau}{k_B}$$

and so

$$(33) \quad R := \frac{k_B}{M}$$

so, as Polk says,  $R$  is different for different gases.

And so

$$\frac{\gamma p}{\rho} = \frac{\gamma \tau}{M}$$

## 11. ISENTROPIC FLOW EQNS. WITH AREA CHANGE

$$\begin{aligned}\frac{p}{p_i} &= \left(\frac{\tau}{\tau_i}\right)^{\frac{\gamma}{\gamma-1}} \\ \frac{\tau}{\tau_i} &= \left(\frac{V_i}{V}\right)^{\gamma-1} = \left(\frac{\rho}{\rho_i}\right)^{\gamma-1} \\ \frac{p}{p_i} &= \left(\frac{V_i}{V}\right)^{\gamma} = \left(\frac{\rho}{\rho_i}\right)^{\gamma}\end{aligned}$$

cf. “Isentropic relations of ideal gas” of Kittel and Kroemer [8]

cf. 20151105 AE121 Polk

$$\begin{aligned}h_0 &= h + \frac{v^2}{2} \\ C_p T_0 &= C_p T + \frac{v^2}{2} \\ \frac{T_0}{T} &= 1 + \frac{v^2}{2C_p T}\end{aligned}$$

Define  $\mathfrak{M} = \frac{v}{a}$ , Mach #

sound speed  $a = (\gamma RT)^{1/2}$

$$(34) \quad \boxed{\frac{T_0}{T} = 1 + \frac{\gamma-1}{2}\mathfrak{M}^2}$$

$$\begin{aligned}\frac{p_0}{p} &= \left(1 + \frac{\gamma-1}{2}\mathfrak{M}^2\right)^{\frac{\gamma}{\gamma-1}} \\ \frac{\rho_0}{\rho} &= \left(1 + \frac{\gamma-1}{2}\mathfrak{M}^2\right)^{\frac{1}{\gamma-1}}\end{aligned}$$

(EY : 20151118 says

$$\begin{aligned}h &:= \frac{H}{MN} \\ c_p &:= \frac{C_p}{MN} \\ H_0 &= H + \frac{MNu^2}{2} \\ H &= C_p \tau\end{aligned} \quad \Longrightarrow \quad \begin{aligned}h_0 &= h + \frac{u^2}{2} \\ c_p \tau_0 &= c_p \tau + \frac{u^2}{2}\end{aligned}$$

and now

$$\begin{aligned}\mathfrak{M} &= \frac{u}{a} = \frac{u}{(\gamma\tau/M)^{1/2}} & C_p &= C_V + N \\ a &= \left(\frac{\gamma\tau}{M}\right)^{1/2} = \left(\frac{\gamma k_B T}{M}\right)^{1/2} = (\gamma RT)^{1/2} & \implies \gamma - 1 &= \frac{N}{C_V}\end{aligned}$$

and so

$$\frac{\tau_0}{\tau} = 1 + \frac{u^2}{2c_p \tau} = 1 + \frac{1}{2c_p \tau} \mathfrak{M}^2 \frac{\gamma\tau}{M} = 1 + \frac{\gamma \mathfrak{M}^2}{2C_p/N} = 1 + \frac{\mathfrak{M}^2}{2C_V/N} = 1 + \frac{\gamma-1}{2} \mathfrak{M}^2$$

) Now we need to get  $\mathfrak{M}$  in terms of area so we can apply these to a nozzle.

Use continuity:

$$\rho_1 v_1 A_1 = \rho_2 A_2 v_2$$

Thus

$$\frac{A_2}{A_1} = \frac{\rho_1 v_1}{\rho_2 v_2} = \frac{\mathfrak{M}_1}{\mathfrak{M}_2} \left(\frac{T_1 \rho_1^2}{T_2 \rho_2^2}\right)^{1/2}$$

EY : 20151120 One can also relate a point in the flow, 1, to another point “downstream” to the flow, 2:

$$\frac{T_1}{T_2} = \frac{T_1/T_0}{T_2/T_0}$$

Substitute for  $T_1/T_2$  and  $p_1/p_e$ ,  $e$  or exh for exhaust,

$$(35) \quad \frac{A_2}{A_1} = \frac{\mathfrak{M}_1}{\mathfrak{M}_2} \left[ \left( \frac{1 + \frac{\gamma-1}{2}\mathfrak{M}_2^2}{1 + \frac{\gamma-1}{2}\mathfrak{M}_1^2} \right)^{\frac{\gamma+1}{\gamma-1}} \right]^{1/2}$$

To get thrust, we need an expression for mass flow rate

$$\begin{aligned}\dot{m} &= \rho v A \\ \frac{\dot{m}}{A} &= \rho v \quad (\text{Recall}) \quad a = (\gamma RT)^{1/2}\end{aligned}$$

$$v = \mathfrak{M}a = \mathfrak{M}(\gamma RT)^{1/2} = \mathfrak{M}(\gamma RT_0)^{1/2} \left(\frac{T}{T_0}\right)^{1/2} = \mathfrak{M} \left[ \frac{\gamma RT_0}{1 + \frac{\gamma-1}{2}\mathfrak{M}^2} \right]^{1/2}$$

Now

$$\begin{aligned}\rho &= \rho_0 \left(\frac{\rho}{\rho_0}\right) = \rho_0 \left[ \frac{1}{1 + \frac{\gamma-1}{2}\mathfrak{M}^2} \right]^{\frac{1}{\gamma-1}} \\ \rho v &= \frac{\dot{m}}{A} = \frac{\rho_0 \mathfrak{M} (\gamma RT_0)^{1/2}}{\left(1 + \frac{\gamma-1}{2}\mathfrak{M}^2\right)^{1/2} \left(1 + \frac{\gamma-1}{2}\mathfrak{M}^2\right)^{\frac{1}{\gamma-1}}}\end{aligned}$$

Using the ideal gas law

$$\frac{\dot{m}}{A} = \frac{p_0 \gamma^{1/2}}{(RT_0)^{1/2}} \mathfrak{M} \frac{1}{\left(1 + \frac{\gamma-1}{2}\mathfrak{M}^2\right)^{\frac{\gamma+1}{2(\gamma-1)}}}$$

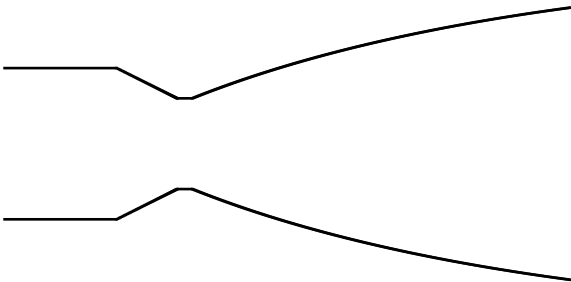
At the throat,  $\mathfrak{M} = 1$ ,

$$(36) \quad \frac{\dot{m}}{A^*} = \frac{p_0 \gamma^{1/2}}{(RT_0)^{1/2}} \left(\frac{2}{\gamma+1}\right)^{\frac{\gamma+1}{2(\gamma-1)}}$$

Thus, from Eq. 36 above (giving the mass flow from throat area and stagnation  $p_0$ ,  $T_0$ ), and plugging this into thrust (force) equation,

$$(37) \quad \begin{aligned}T &= \dot{m} v_e + (p_e - p_a) A_e = \\ &= \frac{A^* p_0 \gamma^{1/2}}{(RT_0)^{1/2}} \left(\frac{2}{\gamma+1}\right)^{\frac{\gamma+1}{2(\gamma-1)}} \sqrt{\frac{2\gamma RT_0}{\gamma-1} \left(1 - \left(\frac{p_e}{p_0}\right)^{\frac{\gamma-1}{\gamma}}\right)} + (p_e - p_a) A_e\end{aligned}$$





12. PSs

12.1. **PS2. Problem 1: Ares I launch vehicle.** Consider a 2-stage rocket. It’s also interesting to explore the properties of polybutadiene acrylonitrile (PBAN) (first stage solid rocket motor fuel), liquid oxygen/liquid hydrogen (LOX/LH2) (second stage liquid rocket engine fuel).

$\epsilon$ , according to wikipedia <sup>4</sup>, is the “ratio between the empty mass of the stage, and the combined empty mass and propellant mass”, which is, in the notation of Biblarz and Sutton (2001) [3],  $\zeta \equiv \frac{m_p}{M(0)}$ .

Neglect drag and earth’s gravity.

Assume constant mass flow.

Let

$$M(t) = M_2(t) + M_{\text{payl}} + M_{01} - m_1(t)$$

with  $M = M(t)$  being the total mass of the entire system during the first stage, with  $m_1 = m_1(t) \in \mathbb{R}$  s.t.

$$m_1(0) = 0$$

$$m_1(t_{p1}) = m_{p1} \text{ (total mass of propellant of stage 1)}$$

$$t_{p1} \equiv \text{burn time of first stage}$$

$$M_{01} \equiv \text{total mass of empty stage 1 + propellant mass for stage 1}$$

$$M_{01} - m_{p1} \equiv \text{mass of empty stage 1.}$$

Also,

$$\epsilon_1 = \frac{M_{01} - m_{p1}}{M_{01}} = 1 - \frac{m_{p1}}{M_{01}}$$

and

$$I_{sp} = \frac{\int F_{\text{thrust}} dt}{W} = \frac{F_{\text{thrust}} t_p}{g_0 \dot{m} t_p} = \frac{u_e}{g_0}$$

Now for the second stage,

$$M_2(t) = \begin{cases} M_{02} & \text{if } 0 \leq t \leq t_{p1} \\ M_{02} - m_2(t) & \end{cases}$$

for  $m_2(t_{p1}) = 0$

$$m_2(t_{p2} + t_{p1}) = m_{p2}$$

<sup>4</sup>[https://en.wikipedia.org/wiki/Multistage\\_rocket](https://en.wikipedia.org/wiki/Multistage_rocket)

From physics, equating kinematics, dynamics  $M\dot{u}$  to the external force,  $F_{\text{thrust}}$ ,

$$M\dot{u} = F_{\text{thrust}} = \dot{m}_1(I_{sp1}g_0) = \frac{m_{p1}}{t_{p1}}I_{sp1}g_0 \text{ so}$$

$$\dot{u} = \frac{\frac{m_{p1}}{t_{p1}}I_{sp1}g_0}{M_2(t) + M_{\text{payl}} + M_{01} - \frac{m_{p1}}{t_{p1}}t} \implies \Delta u_1 = -I_{sp1}g_0 \ln\left(\frac{M_2(t) + M_{\text{payl}} + M_{01} - m_{p1}}{M_2(t) + M_{\text{payl}} + M_{01}}\right)$$

Also,

$$\Delta u_2 = -I_{sp2}g_0 \ln\left(\frac{M_{\text{payl}} + M_{02} - m_{p2}}{M_{\text{payl}} + M_{02}}\right)$$

(a) 

```
gstd = FundConst[ FundConst["Quantity"].str.contains("gravity") ].loc[303,:].Value
# get standard acceleration of gravity
M_0 = Symbol('M_0', positive=True)
Deltau = -I_sp*g_0*ln( (M_0 -m_p)/M_0)
# part (a)
Deltau.subs(I_sp,268.8).subs(g_0,gstd).subs(M_0,805309.).subs(m_p, (1-0.1396)*586344)
# 2595.74521034101 m/s
```

$$\Delta u_1 = 2595.7 \, m/s$$

(b) 

```
Deltau.subs(I_sp,452.1).subs(g_0,gstd).subs(M_0,183952+35013.).subs(m_p, (1-0.1110)*183952)
# 6090.68716730318 m/s
```

$$\Delta u_2 = 6090.7 \, m/s \text{ and so}$$

$$u_{2f} = 8686.4 \, m/s$$

(c) Now

$$\frac{F_{\text{thrust}}}{W} = \frac{\dot{m}_1 u_{e1}}{g_0 M(0)} = 1.5$$

is the thrust to weight ratio at the instant of takeoff.

Now  $u_{e1} = I_{sp1}g_0$  and so

$$\dot{m}_1 = \frac{1.5M(0)}{I_{sp1}}$$

and so

```
>>> 1.5*805309./268.8
# 4493.911830357143
```

so

$$\dot{m}_1 = 4493.9 \, kg/s$$

Over 4 tons of propellant reactants is dumped out per second!

**Problem 2. Continuous staging**

(a)  $M_{\text{payl}} \equiv$  payload mass;  $M_{\text{payl}} \in \mathbb{R}^+$ .

Assume structure mass discarded at 0 velocity.

structure mass continuously jettisoned during the burn.

$M = M(t) \in C^\infty(\mathbb{R})$  represents mass of structure undiscarded at time  $t$ , i.e. system of structure + propellant, at “control volume” at time  $t$ . So consider

$$M = M_{\text{payl}} + m_s(t) + m_p(t) \text{ s.t.}$$

$$m_s = m_s(t) \in C^\infty(\mathbb{R}), \text{ mass of structure not yet thrown out, and propellant mass } m_p = m_p(t) \in C^\infty(\mathbb{R})$$

If we assume constant propellant burn (out), constant propellant flow rate, then

$$\dot{m}_p(t) = \frac{-m_p}{t_p}$$

with  $m_p \in \mathbb{R}^+$  total mass of propellant  
 $t_p \in \mathbb{R}^+$  burn time of propellant fuel.

$$m_p(t) = m_p - \frac{m_p}{t_p}t = m_p(1 - \frac{t}{t_p})$$

Assume constant dead mass ratio  $\delta = m_s(t)/m_p(t)$  (EY: my intuition is we're throwing out as much propellant out in fixed proportion to structure being discarded out; the name "dead" refers to what's still left that's being propelled forward by thrust, I think (???))

$$m_s(t) = \delta m_p(t) = \delta(m_p)(1 - t/t_p)$$

Assume  $I_{sp}$  constant

$$I_{sp} \equiv \frac{\int F_{\text{thrust}} dt}{W_{\text{propellant}}} = \frac{(-\dot{m}_p)u_e t_p}{m_p g_0} = \frac{u_e}{g_0}$$

Consider the instantaneous rest frame of spacecraft + propellant fuel system  $\sum_{\alpha} \Delta p_{i\alpha} = 0$

Consider before and after, after an instant. So for  $\sum_{\alpha} p_{f\alpha}$ ,

$$\text{Now } M = M(t) = M_{\text{payl}} + m_s(t) + m_p(t)$$

$$M(t + \delta t) = M(t) + \dot{M}\delta t$$

Then the momentum of the part that's going to be propelled by the thrust at time  $t + \delta t$  is

$$M(t + \delta t)u(t + \delta t) = (M(t) + \dot{M}\delta t)(u(t) + \dot{u}\delta t) = (M(t) + \dot{M}\delta t)(0 + \dot{u}\delta t) = M\dot{u}\delta t + O((\delta t)^2)$$

The momentum of the propellant expelled out + structure that's discarded is

$$(-\dot{m}_p dt)(-u_e) + (-\dot{m}_s dt) \cdot 0 = -\dot{m}_p dt(-u_e)$$

$$\implies M\dot{u}\delta t + \dot{m}_p u_e \delta t = 0 \implies \dot{u} = -\frac{\dot{m}_p u_e}{M_{\text{payl}} + m_s(t) + m_p(t)} = \frac{\frac{m_p}{t_p} u_e}{M_{\text{payl}} + (\delta + 1)m_p(t)} = \frac{\frac{m_p u_e}{t_p}}{M_{\text{payl}} + (1 + \delta)m_p(1 - \frac{t}{t_p})}$$

$$\implies \Delta u = \frac{-u_e}{1 + \delta} \left[ \ln(M_{\text{payl}} + (1 + \delta)m_p(1 - \frac{t}{t_p})) - \ln(M_{\text{payl}} + (1 + \delta)m_p) \right] = \frac{I_{sp}g_0}{1 + \delta} \ln \left[ \frac{M_{\text{payl}} + (1 + \delta)m_p}{M_{\text{payl}} + (1 + \delta)m_p(1 - \frac{t}{t_p})} \right]$$

$$\text{For } t = t_p, \Delta u = \frac{I_{sp}g_0}{1 + \delta} \ln \left[ \frac{M_{\text{payl}} + (1 + \delta)m_p}{M_{\text{payl}}} \right]$$

(b) Payload fraction is defined as the weight of payload over the takeoff weight (cf. wikipedia <sup>5</sup>). Then for part (a), the  $\Delta u$ ,  $\Delta u_a$ , is

$$(\Delta u)_a = \frac{I_{sp}g_0}{1 + \delta} \ln \left[ \frac{1}{\lambda} \right] = \frac{-I_{sp}g_0}{1 + \delta} \ln \lambda$$

(i)  $M = M(t) = M_{\text{payl}} + M_s + m_p(t)$  s.t.  $m_p(0) = M_p$  if assume constant mass flow out.  $m_p(t) = M_p(1 - \frac{t}{t_p})$   
 $m_p(t_p) = 0$

$$M\dot{u} = -\dot{m}_p u_{\text{exh}} \implies \dot{u} = \frac{\frac{M_p}{t_p} u_e}{M_{\text{payl}} + M_s + M_p(1 - \frac{t}{t_p})}$$

So

$$u(t) = -I_{sp}g_0 [\ln(M_{\text{payl}} + M_s + M_p(1 - \frac{t}{t_p})) - \ln(M_{\text{payl}} + M_s + M_p)]$$

$$\Delta u = I_{sp}g_0 \ln \left[ \frac{M_{\text{payl}} + M_s + M_p}{M_{\text{payl}} + M_s} \right]$$

Then it's just algebra to put the ratio of the masses above in terms of  $\lambda, \delta$  (there's 3 unknowns,  $M_{\text{payl}}$ ,  $M_s$ ,  $M_p$  masses, and we're given  $\lambda, \delta$  and a ratio we want,  $\frac{M_{\text{payl}} + M_s + M_p}{M_{\text{payl}} + M_s}$ ). Instead of doing the algebra entirely by hand, let's use Python and the sympy library:

```
import sympy
from sympy import *
>>> M_payl = Symbol('M_payl', positive=True)
>>> M_s = Symbol('M_s', positive=True)
>>> ratio_bi = (M_payl + M_s)/(M_payl + M_s + M_p)
>>> payloadfrac = Symbol('payloadfrac', positive=True)
>>> delta = Symbol('delta', positive=True)
>>> ratio_bi_new = ratio_bi.subs(M_payl, (M_s + M_p)/(1/payloadfrac - 1)).subs(M_s, M_p*delta)
>>> ratio_bi_new.expand().factor(M_p).simplify().factor()
(delta + payloadfrac)/(delta + 1)
```

and so

$$(\Delta u)_i = I_{sp}g_0 \ln \left[ \frac{1 + \delta}{\lambda + \delta} \right]$$

(ii)

$$(\Delta u)_{ii} = I_{sp}g_0 \ln \left[ \frac{M_{\text{payl}} + M_p}{M_{\text{payl}}} \right] = I_{sp}g_0 \ln \left( \frac{1}{\lambda} \right)$$

(iii) For 1 (total number of) stage(s),

$$\Delta u = I_{sp}g_0 \ln \left[ \frac{1 + \delta}{\lambda + \delta} \right]$$

as in part (b)(i).

Let  $N$  be the number of stages.

**12.2. Lagrangian point of view for gravity-free, drag-free rocket.** Let  $M = M(t) = M(0) - m(t)$  s.t.  $m(0) = 0$  where  $m(t_p) = M_p$

$M_p$  is the total mass of the propellant to be used, and  $M(0)$  represents the (initial) mass of the propellant + spacecraft or payload or the interesting part we want to launch out.  $t_p$  is the burn time of the propellant.

Then kinetic energy of  $M$  at instantaneous time  $t$  is  $\frac{1}{2}Mu^2$ . Also note that  $\int_0^{t_p} \dot{m} dt = M_p$ .

Now take the exterior derivative of  $m$ :  $dm = \dot{m}dt$ . Consider the infinitesimal piece of mass  $dm$  at the (instantaneous) time  $t$ ; its kinetic energy will be

$$\frac{1}{2}dm(u - u_e)^2$$

Note that its velocity is  $(u - u_e)$  because mass  $dm$  is being expunged out of the rocket at constant exhaust velocity  $u_e$  *relative* to the rocket (i.e. change, transform, or "boost" into the instantaneous inertial reference frame of the rocket, so that the rocket has 0 velocity in this frame; now return back to the "lab" frame-the propellant expunged has gained velocity  $u$  as well).

Now

$$\frac{1}{2}dm(u - u_e)^2 \xrightarrow{\int_0^t dt} \frac{1}{2} \int_0^t dm(u - u_e)^2 = \frac{1}{2} \int_0^t \dot{m}d\tau(u - u_e)^2$$

The full Lagrangian at time  $t$  is

$$\mathcal{L} = \frac{1}{2}Mu^2 + \frac{1}{2} \int_0^t \dot{m}d\tau(u - u_e)^2$$

Notice that  $\mathcal{L}$  in this specific case does not depend on position. So  $\frac{\partial \mathcal{L}}{\partial x^i} = 0$ .

Now

$$\frac{\partial \frac{1}{2}Mu^2}{\partial u^i} = Mu_i \xrightarrow{\frac{d}{dt}} \frac{d}{dt} \frac{\partial \frac{1}{2}Mu^2}{\partial u^i} = \dot{M}u_i + M\dot{u}_i = -\dot{m}u_i + M\dot{u}_i$$

<sup>5</sup>"Payload fraction", *Wikipedia* [https://en.wikipedia.org/wiki/Payload\\_fraction](https://en.wikipedia.org/wiki/Payload_fraction)

and

$$\frac{\partial}{\partial u^i} \frac{1}{2} \int_0^t \dot{m} dt (u - u_e)^2 = \frac{1}{2} \int_0^t \dot{m} d\tau 2(u_i - u_e) = \int_0^t \dot{m} d\tau (u_i - u_e)$$
$$\xrightarrow{\frac{d}{dt}} \frac{d}{dt} \frac{\partial}{\partial u^i} \frac{1}{2} \int_0^t \dot{m} dt (u - u_e)^2 = \frac{d}{dt} \int_0^t \dot{m} d\tau (u_i - u_e) = \dot{m}(t)(u_i - u_e)$$

for clearly, if  $\int_0^t \dot{m} d\tau (u_i - u_e) = F(t) - F(0)$ , then applying the total time derivative, will result in  $F'(t)$  and so we can read off the antiderivative, namely  $\dot{m}(t)(u_i - u_e)$  at time  $t$ .

The Euler-Lagrange equation tells us that  $\frac{\partial \mathcal{L}}{\partial x^i} - \frac{d}{dt} \frac{\partial \mathcal{L}}{\partial \dot{x}^i} = 0$  and so

$$-\dot{m} u_i + M \dot{u}_i + \dot{m} u_i - \dot{m} u_e = 0 \text{ or } M \dot{u}_i = \dot{m} u_e$$

$$\implies \boxed{\dot{u}_i = \frac{\dot{m} u_e}{M}}$$

Problem 4. Ballistic trajectories with atmospheric drag

(a) Now

$$M_0 \dot{u} = -M_0 g + F_d$$

$$\dot{u} = -g + \frac{F_d}{M_0}$$

In components,

$$\dot{u}_x = \frac{F_d}{M_0} \left( \frac{-u_x}{|u|} \right)$$

$$\dot{u}_y = -g + \frac{F_d}{M_0} \left( \frac{-u_y}{|u|} \right)$$

And so, for

$$F_d = \frac{1}{2} \rho C_D u^2 A = \frac{1}{2} \rho C_D (u_x^2 + u_y^2) A$$

and

$$\rho = \rho_0 e^{-y/\lambda}$$

$$\ddot{x} = \frac{1}{2M_0} \rho_0 e^{-y/\lambda} C_D (\dot{x}^2 + \dot{y}^2) A \left( \frac{-\dot{x}}{(\dot{x}^2 + \dot{y}^2)^{1/2}} \right)$$

$$\ddot{y} = -g + \frac{1}{2M_0} \rho_0 e^{-y/\lambda} C_D (\dot{x}^2 + \dot{y}^2) A \left( \frac{-\dot{y}}{(\dot{x}^2 + \dot{y}^2)^{1/2}} \right)$$

Note that we can rewrite this as the following system of equations:

$$\dot{u}_x = \frac{1}{2M_0} \rho_0 C_D A e^{-y/\lambda} (u_x^2 + u_y^2)^{1/2} (-u_x)$$

$$\dot{u}_y = -g_0 + \frac{1}{2M_0} \rho_0 C_D A e^{-y/\lambda} (u_x^2 + u_y^2)^{1/2} (-u_y)$$

$$\dot{x} = u_x$$

$$\dot{y} = u_y$$

In Propulsion.py, (just run it with `python -i Propulsion.py` in its directory)

```
import scipy
from scipy import exp, array
from scipy.integrate import ode
```

```
import matplotlib.pyplot as plt
```

```
M_cannonball = (7.8*(10**2)**3/(10**3))*4./3.*N(pi)*(15./2./100. )**3
(1.225)*(0.1)/(2.*M_cannonball)*(N(pi)*(15./2./100. )**2) # 7.85256410256411e-5
```

```
# to use scipy.integrate.ode, the ODE system must be declared as a "callable", a Python function
def deriv(t,u): # return derivatives of the array u
    """
    cf. http://bulldog2.redlands.edu/facultyfolder/deweerd/tutorials/Tutorial-ODEs.pdf

    """
    uxdot = (7.853*10**(-5))*exp( -u[3]/(10000.))*(u[0]**2 + u[1]**2)**(0.5)*(-u[0])
    uydots = -9.8 + (7.853*10**(-5))*exp(-u[3]/(10000.))*(u[0]**2 + u[1]**2)**(0.5)*(-u[1])
    return array([ uxdot, uydots, u[0], u[1] ])
```

```
# initial conditions
u0 = [300.*cos(50./180.*N(pi)), 300.*sin(50./180.*N(pi)),0,0]

# declare the ODE to be integrated
Prob0203 = ode(deriv).set_integrator('dopri5') # Problem 3 from Problem Set 2 for AE121 Fall 2015
# cf. http://stackoverflow.com/questions/26738676/does-scipy-integrate-ode-set-solout-work
Prob0203.set_initial_value(u0)
```

```
t1 = 41.575
dt = 0.005
# print out the solution to the ODE for various times as a sanity check
while Prob0203.successful() and Prob0203.t < t1:
    Prob0203.integrate(Prob0203.t+dt)
    print("%g-" % Prob0203.t )
    print Prob0203.y
```

```
# store the solutions in a Python list for plotting
Prob0203.set_initial_value(u0)
Prob0203.solution = []
while Prob0203.successful() and Prob0203.t < t1:
    Prob0203.solution.append( [Prob0203.t+dt,] + list( Prob0203.integrate(Prob0203.t+dt) ) )
# take the transpose of a list of lists
Prob0203.solution = map(list, zip(*Prob0203.solution))
```

```
plt.figure(1)
plt.plot( Prob0203.solution[3], Prob0203.solution[4])
plt.xlabel('x-(m)')
plt.ylabel('y-(m)')
plt.title(' Cannonball_trajectory_with_Drag: Variable_density')
```

Horizontal range is 6.11 km

(b)

```
def deriv_b(t,u): # return derivatives of the array u
    """
    cf. http://bulldog2.redlands.edu/facultyfolder/deweerd/tutorials/Tutorial-ODEs.pdf

    """
    uxdot = (7.853*10**(-5)) *(u[0]**2 + u[1]**2)**(0.5)*(-u[0])
    uydots = -9.8 + (7.853*10**(-5)) *(u[0]**2 + u[1]**2)**(0.5)*(-u[1])
    return array([ uxdot, uydots, u[0], u[1] ])
```

```
Prob0203b = ode(deriv_b).set_integrator('dopri5')
Prob0203b.set_initial_value(u0)
Prob0203b.integrate(41.23)
```

```
t1b = 41.225
Prob0203b.set_initial_value(u0)
Prob0203b.solution = []
while Prob0203b.successful() and Prob0203b.t < t1b:
    Prob0203b.solution.append( [Prob0203b.t+dt,] + list( Prob0203b.integrate(Prob0203b.t+dt) ) )
Prob0203b.solution = map(list, zip(*Prob0203b.solution))
```

```
plt.figure(2)
plt.plot( Prob0203b.solution[3], Prob0203b.solution[4])
plt.xlabel('x-(m)')
plt.ylabel('y-(m)')
```

```
plt.title('Cannonball_trajectory_with_Drag:_Constant_density')
```

Horizontal range is  $5.89\text{ km}$ . This makes sense because the cannonball finds it easier to fly “through the air” at higher altitudes, higher up the atmosphere, because the “air is thinner” in the “upper atmosphere.”

(c) For no drag, this can be solved analytically:

$$\begin{aligned} \dot{u}_x &= 0 & u_x &= u_0 \cos \theta & x(t) &= u_0 \cos \theta t \\ \dot{u}_y &= -g & \implies u_y &= -gt + u_0 \sin \theta & \implies y(t) &= -\frac{1}{2}gt^2 + u_0 \sin \theta t = t(u_0 \sin \theta - \frac{gt}{2}) \end{aligned}$$

So the horizontal range is  $x(t_f) = u_0 \cos \theta \left( \frac{2u_0 \sin \theta}{g} \right) = \frac{u_0^2}{g} \sin(2\theta) = 9044.m$  for

```
>>> 300.**2/9.8*sin(2.*50./180.*N(pi))
9044.15283378558
```

```
#parabola trajectory data
Prob0203c_x = [i*10 for i in range(905)]
Prob0203c_y = [tan(50./180.*N(pi))*x - (9.8/2.)*x**2/(300.*cos(50./180.*N(pi)))*2 for x in Prob0203c_x]

# plot all 3 trajectories together
plt.figure(3)
plt.plot(Prob0203_solution[3], Prob0203_solution[4], label="Drag:_Variable_density")
plt.plot(Prob0203b_solution[3], Prob0203b_solution[4], label="Drag:_Constant_density")
plt.plot(Prob0203c_x, Prob0203c_y, label="No_Drag")
plt.xlabel('x_(m)')
plt.ylabel('y_(m)')
plt.title('Trajectories_of_cannonball_with_Drag_of_variable_density,_Drag_of_constant_density,_and_no_drag')
plt.legend()
```

If there was no drag, then the cannonball will fly out farther, and higher. It’s important to consider air resistance, as the horizontal range difference between drag and no drag is almost 3000 m (!!!). It’s important to consider variation of atmospheric drag with altitude for horizontal range for precision landing (about a 300 m difference).

12.3. PS4. Problem 1: Kinetic theory connection to thermodynamics properties.

- (a)  
(b) Now the Maxwellian velocity distribution,  $P(v)$ , where  $P(v)dv$  is the probability that the particle has speed in  $(v, v+dv)$ , is given by

$$P(v) = 4\pi \left( \frac{M}{2\pi\tau} \right)^{3/2} v^2 \exp \left( \frac{-Mv^2}{2\tau} \right)$$

If  $N$  particles had the same kinetic energy, then the entire system of  $N$  particles would have a total internal energy of  $N\frac{1}{2}Mv^2$ , with  $M$  being the mass of a single particle.

Thus, the total internal energy  $U$  is calculated by weighting by  $P(v)$  the above total kinetic energy, which in this case of only 3 translational degrees of freedom, coincides with the total internal energy:

$$\begin{aligned} U &= N \int_0^\infty dv \left( \frac{1}{2}Mv^2 \right) P(v) = \int_0^\infty dv \frac{MN}{2} \cdot 4\pi \left( \frac{M}{2\pi\tau} \right)^{3/2} v^4 \exp \left( \frac{-Mv^2}{2\tau} \right) = 2\pi MN \left( \frac{M}{2\pi\tau} \right)^{3/2} \int_0^\infty dv v^4 \exp \left( \frac{-Mv^2}{2\tau} \right) = \\ &= 2\pi MN \left( \frac{M}{2\pi\tau} \right)^{3/2} \frac{3}{4 \left( \frac{M}{2\tau} \right)^2} \frac{\sqrt{\pi}}{2 \left( \frac{M}{2\tau} \right)^{1/2}} = 2\pi MN \frac{3}{4} \frac{1}{\left( \frac{M}{2\tau} \right)} \frac{\sqrt{\pi}}{2\pi^{3/2}} = \frac{3\tau N}{2} \end{aligned}$$

Also,  $\frac{U}{N} = \frac{3\tau}{2}$ .

EY : 20151117 I want to reiterate that there must be a more systematic and sane and rational way of looking up Physical Constants and other physical data than by manually looking it up a book or manually looking it up a website. People can make a mistake copying and pasting! Thus, I wrote the Python package `Physique` that you copy into your working directory and can import in (this is all in the script `Propulsion.py`:

```
import Physique
from Physique import FCconv, KCconv, FundConst, conv, plnfacts, T_C, T_K, T_F

k_Boltz = FundConst[ FundConst["Quantity"].str.contains("Boltzmann") ].loc[49,: ]
>>> k_Boltz.Value
Decimal('1.38064852E-23')
>>> k_Boltz.Unit
'J_K^-1'
```

So for  $T = 300K$  and  $T = 1000K$ ,

```
>>> k_Boltz.Value *300*Decimal(1.5)
Decimal('6.21291834000E-21')
>>> k_Boltz.Value *1000*Decimal(1.5)
Decimal('2.070972780000E-20')
```

or

```
>>> k_Boltz.Value *300*Decimal(1.5)/ JovereV.Multiplyby
Decimal('0.03877797733958233079116726804')
>>> k_Boltz.Value *1000*Decimal(1.5)/ JovereV.Multiplyby
Decimal('0.1292599244652744359705575601')
```

and so  $U/N = 6.213 \times 10^{-21}J$  or  $0.0388eV$  for  $T = 300K$  and  $2.071 \times 10^{-20}J$  or  $0.129eV$  for  $T = 1000K$

(c) Now

$$C_V = \left( \frac{\partial U}{\partial \tau} \right)_V = \frac{3N}{2}$$

From the definition of  $C_V$ . Also, for enthalpy  $H$ ,  $H = U + pV = U + N\tau$ , for ideal gas law still holds,

$$C_P = \left( \frac{\partial H}{\partial \tau} \right)_V = \left( \frac{\partial U}{\partial \tau} \right)_V + N = C_V + N$$

and so

$$C_P = \frac{5N}{2}$$

These are all, above, physicists’ quantities. For engineers, *specific heat capacities* are useful with real-world material. *Specific* heat capacities are obtained from physicists’ heat capacities by dividing by what physicists would’ve deemed as  $M$ , the (straight-up) mass.

$$\begin{aligned} c_V &:= \frac{C_V}{M} = \frac{3N}{2M} \\ c_P &:= \frac{C_P}{M} = \frac{5N}{2M} \end{aligned}$$

```
N_Avog = FundConst[FundConst['Quantity'].str.contains('Avogadro')]
>>> c_V = float(Decimal(1.5)*(N_Avog.Value)*(k_Boltz.Value))/M_0
>>> c_P = float(Decimal(2.5)*(N_Avog.Value)*(k_Boltz.Value))/M_0
>>> c_V.subs(M_0, 39.948/1000.)
312.198102337360
>>> c_V.subs(M_0, 131.293/1000.)
94.9912774647001
>>> c_P.subs(M_0, 39.948/1000.)
520.330170562267
>>> c_P.subs(M_0, 131.293/1000.)
158.318795774500
```

So  $c_V$  and  $c_P$  for argon is  $312.2J/(kgK)$  and  $520.3J/(kgK)$ , respectively, and  $c_V$  and  $c_P$  for xenon is  $94.99J/(kgK)$  and  $158.3J/(kgK)$ , respectively.

The so-called *molar heat capacity* is the amount of heat needed to increase the temperature of 1 mole of substance. Now, physicists’  $C_V$  and  $C_P$  is the heat capacities for the amount of heat needed to raise the temperature of a system of  $N$  number

of particles. Then certainly, dividing  $C_V, C_P$  by  $N$  will result in the amount of heat needed to raise the temperature of a *single* particle. Use Avogadro’s number to convert between number of particles and moles.

$$c_V = \frac{C_V}{N} = \frac{3}{2}$$
$$c_P = \frac{C_P}{N} = \frac{5}{2}$$

```
>>> Decimal(1.5)*(N_Avog.Value)*(k_Boltz.Value)
42      12.471689792172872460
Name: Value, dtype: object
>>> Decimal(2.5)*(N_Avog.Value)*(k_Boltz.Value)
42      20.786149653621454100
```

and so  $c_V = 12472.J/(\text{kmol}K)$  and  $c_P = 20786.J/(\text{kmol}K)$

**Problem 2: Mean thermal velocity.**

Wikipedia article “Thermal velocity” has 3 mean thermal velocities <sup>6</sup>

- $P(v) = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} v^2 \exp\left(\frac{-Mv^2}{2\tau}\right)$   
Now

$$\int_0^\infty dv v^4 \exp(-\alpha v^2) = \int_0^\infty v^3 \left(\frac{\exp(-\alpha v^2)}{-2\alpha}\right)' = 0 - \int 3v^2 \frac{\exp(-\alpha v^2)}{-2\alpha} = \int_0^\infty \frac{3}{2\alpha} v^2 \exp(-\alpha v^2) =$$
$$= \frac{3}{2\alpha} \int v \left(\frac{e^{-\alpha v^2}}{-2\alpha}\right)' = \frac{3}{2\alpha} \left[0 - \int \frac{e^{-\alpha v^2}}{-2\alpha}\right] = \frac{3}{4\alpha^2} \frac{\sqrt{\pi}}{2\sqrt{\alpha}}$$

So

$$\langle v^2 \rangle = \int_0^\infty v^2 P(v) = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \frac{3}{4 \left(\frac{M}{2\tau}\right)^2} \frac{\sqrt{\pi}}{\sqrt{\frac{M}{2\tau}}} \frac{1}{2} = \frac{3}{\frac{M}{2\tau}} \frac{1}{2} = \frac{3\tau}{M}$$

so

$$v_{rms} = \left(\frac{3\tau}{M}\right)^{1/2}$$

Now

$$\langle v^2 \rangle = 3\langle v_x^2 \rangle = \frac{3\tau}{M} \text{ so } \langle v_x^2 \rangle^{1/2} = \left(\frac{\tau}{M}\right)^{1/2}$$

- 

$$\frac{\partial P(v)}{\partial v} = \frac{2P(v)}{v} + P(v) \left(\frac{-M}{\tau}\right) v = P(v) \left[\frac{2}{v} - \frac{M}{\tau} v\right] = 0 \text{ if } v_{mp} = 0 \text{ or } v_{mp} = \sqrt{\frac{2\tau}{M}}$$

where  $mp$  stands for most probable.

$$v_{mp} = \sqrt{\frac{2\tau}{M}} < \sqrt{\frac{3\tau}{M}} = v_{rms}$$

<sup>6</sup>“Thermal Velocity”, *wikipedia*, [https://en.wikipedia.org/wiki/Thermal\\_velocity](https://en.wikipedia.org/wiki/Thermal_velocity)

- 

$$\bar{c} = \int_0^\infty dv v P(v) =: \langle v \rangle = \int_0^\infty dv 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} v^3 \exp\left(\frac{-Mv^2}{2\tau}\right) = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \int_0^\infty dv v^2 \left(\frac{e^{-\frac{Mv^2}{2\tau}}}{-\frac{M}{\tau}}\right)' =$$
$$= 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \left[0 - \int_0^\infty dv 2v \frac{e^{-\frac{Mv^2}{2\tau}}}{-\frac{M}{\tau}}\right] = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \left[\frac{2\tau}{M} \left(\frac{e^{-\frac{Mv^2}{2\tau}}}{-\frac{M}{\tau}}\right)\right]_0^\infty = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \frac{2\tau^2}{M^2} =$$
$$= \frac{1}{(2\pi)^{3/2}} 8\pi \sqrt{\frac{\tau}{M}} = \sqrt{\frac{8\tau}{\pi M}}$$

Note

$$\frac{v_{rms}}{\bar{c}} = \frac{\left(\frac{3\tau}{M}\right)^{1/2}}{\left(\frac{8\tau}{\pi M}\right)^{1/2}} = \left(\frac{3\pi}{8}\right)^{1/2}$$

**Problem 3: Ideal vs. real rocket analysis.**

**12.4. PS5. Problem 1: Nozzle flow in liquid rocket engines.**

Viking series liquiad rocket engines used on first 2 stages of Ariane 4 launch vehicle.  
Rocket engines are storable propellant motors that use  
nitrogen tetroxide and UDMH25 (unsymmetrical dimethyl hydrazine with 25 percent hydrazine hydrate) as propellants.  
suitable mean value for molecular weight of the combustion product mixture 23 *g*/mol ratio of specific heats approximately 1.2

Combustion chamber temperatures for both 3350*K* approximately.

Viking 5C motor: chamber pressure 5800 *kPa*

propellant mass flow rate 275.2 *kg/s*

nozzle expansion ratio 10.5

Viking 4B engine: chamber pressure 5850 *kPa*

mass flow 278.0 *kg/s*

expansion ratio 30.8

- (a) Use Eq. **36**,

$$\frac{\dot{m}}{A^*} = \frac{p_0 \gamma^{1/2}}{(RT_0)^{1/2}} \left(\frac{2}{\gamma + 1}\right)^{\frac{\gamma+1}{2(\gamma-1)}}$$

Opening up NozzleTheory.py, using Python’s **sympy** library,

```
Viking5Cnozzle = massflowrateExp.subs(gamma, 1.2).subs(massflow, 275.2).subs(p_0,5800*1000).subs(T_0,3350).\
subs(R, k_Boltz.Value/ (Decimal(23/1000.)/N_Avog.Value ))
Viking4Bnozzle = massflowrateExp.subs(gamma, 1.2).subs(massflow, 278.0).subs(p_0,5850*1000).subs(T_0,3350).\
subs(R, k_Boltz.Value/ (Decimal(23/1000.)/N_Avog.Value ))
```

```
solve(Viking5Cnozzle, Astar)[0] # 0.0805128291046479
solve(Viking4Bnozzle, Astar)[0] # 0.0806368550290290
```

So the area of throat for Viking 5C and Viking 4B is 0.08051 *m*<sup>2</sup> and 0.08126 *m*<sup>2</sup>, respectively. Then, the expansion ratio gives the expanded area for Viking 5C and Viking 4B of 0.8454 *m*<sup>2</sup> and 2.484 *m*<sup>2</sup>, respectively.

- (b) We need a relationship relating (cross-sectional) area ratio to Mach number, Eq. **35**

$$\frac{A_2}{A_1} = \frac{\mathfrak{M}_1}{\mathfrak{M}_2} \left[ \left(\frac{1 + \frac{\gamma-1}{2} \mathfrak{M}_2^2}{1 + \frac{\gamma-1}{2} \mathfrak{M}_1^2}\right)^{\frac{\gamma+1}{\gamma-1}} \right]^{1/2}$$



At the throat,  $\mathfrak{M} = 1$ , and denote the area by  $A^*$ . Thus

$$\frac{A_e}{A^*} = \frac{1}{\mathfrak{M}} \left[ \left( \frac{1 + \frac{\gamma-1}{2} \mathfrak{M}^2}{\frac{\gamma+1}{2}} \right)^{\frac{\gamma+1}{\gamma-1}} \right]^{1/2}$$

and so  $\frac{A_e}{A^*}$  is the expansion ratio.

This is implemented in `NozzleTheory.py` as a sympy object Eq:

```
# Area Ratio to Mach numbers
Mach_1 = Symbol('Mach_1', positive=True)
Mach_2 = Symbol('Mach_2', positive=True)
A_1     = Symbol('A_1',  positive=True)
A_2     = Symbol('A_2',  positive=True)
AreastoMachs = Eq( A_2/A_1 , Mach_1/Mach_2*( ( (Rat(1) + (gamma-Rat(1))/(Rat(2) )*Mach_2**2 )/\
(Rat(1) + (gamma-Rat(1))/(Rat(2) )*Mach_1**2 ) )**(( gamma +1)/(gamma-1) ) )**0.5 )
```

and one substitutes in the desired, given numbers (parameters):

```
Viking5CMachEq = AreastoMachs.subs(Mach_1,Rat(1) ).subs(gamma,1.2).subs(A_2,A_1*10.5)
# 10.5 == (0.0909090909090909*Mach_2**2 + 0.909090909090909)*5.5/Mach_2
```

Notice that now we have a so-called “root-finding” problem, with non-integer exponents. One should use a numerical solver so that finding this root is done “efficiently”. Also, to do this in Python’s `scipy`, we have to create a Python function object, and so I used sympy’s `lambdify` to turn a sympy expression into an actual Python function:

```
Viking5CMach = lambdify(Mach_2, Viking5CMachEq.rhs - Viking5CMachEq.lhs )
# Remember to move all terms to 1 side, and so the other side equals 0
```

```
scipy.optimize.newton( Viking5CMach, 3) # Newton Raphson method
# 3.3123573073570207
```

```
scipy.optimize.bisect( Viking5CMach, 3,4) # Bisection method
# 3.312357307356251
```

[AIMS Senegal](#) had a nice, easy introduction to scipy’s root finding methods <sup>7</sup>.

Next, the static temperature, static pressure, static density at the nozzle exit can be easily calculated from the isentropic relations:

$$\frac{\tau}{\tau_0} = \left( 1 + \frac{\gamma-1}{2} \mathfrak{M}^2 \right)^{-1} \quad \frac{p}{p_0} = \left( \frac{\tau}{\tau_0} \right)^{\frac{\gamma}{\gamma-1}} \quad \rho_0 = \frac{p_0}{RT_0} \quad \frac{\rho}{\rho_0} = \left( \frac{\tau}{\tau_0} \right)^{\frac{1}{\gamma-1}}$$

Then the static pressure, static temperature, static density for Viking 5C and Viking 4B engines are

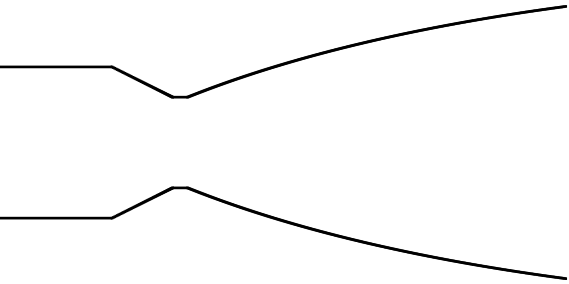
$$\begin{aligned} \mathfrak{M}_{exh} &= 3.312 \\ p &= 68.2 \text{ kPa} \\ T &= 1597.4 \text{ K} \\ \rho &= 0.1181 \text{ kg}/m^3 \\ \mathfrak{M}_{exh} &= 4.057 \\ p &= 17.0 \text{ kPa} \\ T &= 1265.9 \text{ K} \\ \rho &= 0.0372 \text{ kg}/m^3 \end{aligned}$$

(c)

<sup>7</sup>Roots finding, Numerical integrations and differential equations, *AIMS Senegal*, [AIMS Senegal](#)

(d) Look at Eq. [37](#) again:

$$\begin{aligned} T &= \dot{m} v_e + (p_e - p_a) A_e = \\ &= \dot{m} \sqrt{\frac{2\gamma RT_0}{\gamma-1} \left( 1 - \left( \frac{p_e}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right)} + (p_e - p_a) A_e \end{aligned}$$



So if  $p_e = p_a$ , then

$$T = \dot{m} \sqrt{\frac{2\gamma RT_0}{\gamma-1} \left( 1 - \left( \frac{p_a}{p_0} \right)^{\frac{\gamma-1}{\gamma}} \right)}$$

(e) Look at this webpage: <http://www.engapplets.vt.edu/fluids/CDnozzle/cdinfo.html> for a good recap of the physics of shocks, accompanying the lecture by Polk in AE121.

Also “Choked flow” wikipedia article [https://en.wikipedia.org/wiki/Choked\\_flow](https://en.wikipedia.org/wiki/Choked_flow)

Consider normal shock at the exit of the nozzle. Then before,  $\mathfrak{M}_1 > 1$  (supersonic), and right after the shock,  $\mathfrak{M}_2 < 1$  (subsonic)

**Problem 3: Duct flow with heating.**

(a) In this case, I think that the correct heat capacity to use is  $C_p$  because the process occurs at the constant ambient pressure (that ambient pressure, external to the (constant area) duct, remains constant and on the duct).

Now

$$Q = C_p d\tau$$

For  $c_p := \frac{C_p}{MN}$ ,  $q = mc_p dT$  (experimental physicists’ and engineers’ expression). Thus

$$\dot{q} = \dot{m} c_p dT \implies \int_{\gamma} \dot{q} = \Delta \dot{q} = \dot{m} c_p (T_f - T_i) \text{ or } \frac{\Delta \dot{q}}{\dot{m} c_p} = T_f - T_i$$

Now this heating raises the temperature of the *stagnation* (enthalpy) temperature, because we’re considering this heating as a (adiabatic?) thermodynamic process, separate from what’s going on with the flow.

For instance, along the flow, at a point, the thermodynamic property is  $(T, p)$ . We can always “pull” this “back” to the stagnation properties, whether hypothetical or real:

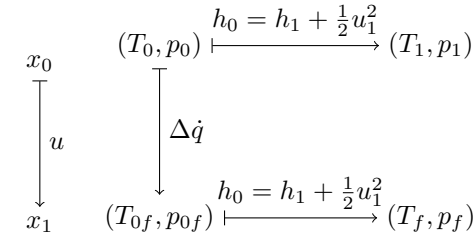


upstream  $\longrightarrow$  downstream

$$(T_0, p_0) \longmapsto (T, p)$$

$$\frac{T_0}{T} = 1 + \frac{\gamma - 1}{2} \mathfrak{M}^2$$
$$\frac{p_0}{p} = \left(1 + \frac{\gamma - 1}{2} \mathfrak{M}^2\right)^{\frac{\gamma}{\gamma - 1}}$$

Thus, when one considers heat addition, the thermodynamic process,  $\Delta \dot{q}$  occurs on stagnation temperature, and then can be related to the actual, physical fluid flow and its properties, through formulae we’ve derived:



```
(1.*10**6)/(40.*15.6)
# 1602.5641025641025
```

So

$$\Delta T_0 = 1602.6\text{ K}$$

- (b) Now consider the converging section at the inlet to each 1 cm diameter channel, that “accelerates the flow to a relatively low Mach number and produces a flow rate of 40 g/s”.

Mass continuity (conservation) still holds:

$$\dot{m} = \rho A u = \frac{p}{RT} A \mathfrak{M} \sqrt{\gamma R T} = \frac{p}{\sqrt{RT}} A \mathfrak{M} \sqrt{\gamma} = \frac{p_0 \sqrt{\gamma}}{\sqrt{RT_0}} A \mathfrak{M} \left(1 + \frac{\gamma - 1}{2} \mathfrak{M}^2\right)^{-\frac{\gamma + 1}{2(\gamma - 1)}}$$

where the “pullback” to the stagnation properties for each point of the flow, before the converging section and after the converging section, was used:

$$\frac{p_0}{p} = \left(1 + \frac{\gamma - 1}{2} \mathfrak{M}^2\right)^{\frac{\gamma}{\gamma - 1}}$$
$$\frac{T_0}{T} = 1 + \frac{\gamma - 1}{2} \mathfrak{M}^2$$

```
massconsEq = Eq(massflow , p_0*sqrt(gamma/(R*T_0))*A*Mach*(1+(gamma-Rat(1))/Rat(2)*Mach**2)**\
((gamma+1)/(-Rat(2)*(gamma-1))))
massconsProb0503 = massconsEq.subs(massflow , 40.*10**(-3)).subs(gamma, 1.4).subs(p_0 , 6.8*10**(-6)).\
subs(T_0 , 673.).subs(A,N(pi)*(10**(-2)/2.))**2).subs(R, k_Boltz.Value/(Decimal(2.0159*10**(-3))/N_Avog.Value))
```

where for  $R$ ,  $R = \frac{k_B}{M}$  and where for  $M$ , I used 2.0159g/mol for H<sub>2</sub>, and Avogadro’s number,  $N_A = 6.022140857 \times 10^{23}$ , which is number of particles per mole. So in this example,  $R = 4124.4$  for H<sub>2</sub> as

```
>>> k_Boltz.Value/(Decimal(2.0159*10**(-3))/N_Avog.Value)
Decimal('4124.440627733807619868515695')
```

Let’s try to use the derived relation, relating the stagnation temperatures before and after (denoted 1) heat addition:

$$\frac{T_0}{T_{01}} = \left[ \frac{1 + \gamma \mathfrak{M}_1^2}{1 + \gamma \mathfrak{M}^2} \left( \frac{\mathfrak{M}}{\mathfrak{M}_1} \right) \right]^2 \left[ \frac{1 + \frac{\gamma - 1}{2} \mathfrak{M}^2}{1 + \frac{\gamma - 1}{2} \mathfrak{M}_1^2} \right]$$

This is implemented in `NozzleTheory.py`:

```
T_01=Symbol('T_01',real=True)
heataddTvsMachEq=Eq(T_0/T_01,((Rat(1)+gamma*Mach_1**2)/(Rat(1)+gamma*Mach**2)*Mach/Mach_1)**2*\
(Rat(1)+(gamma-Rat(1))/Rat(2)*Mach**2)/(1+(gamma-Rat(1))/Rat(2)*Mach_1**2))

heataddTvsMachProb0503=heataddTvsMachEq.subs(gamma,1.4).subs(T_0,675).subs(T_01,675+(1.*10**6)/\
(40.*15.6)).subs(Mach,MachProb0503)
Mach1Prob0503lamb=lambdify(Mach_1,heataddTvsMachProb0503.rhs-heataddTvsMachProb0503.lhs)
plot(heataddTvsMachProb0503.rhs,(Mach_1,0,10))
Mach1Prob0503=scipy.optimize.newton(Mach1Prob0503lamb,0.1)
```

where  $\mathfrak{M}$  was obtained for the acceleration by the converging nozzle,  $T_0$  is the stagnation temperature that was given, and  $T_{01}$  is obtained from part (a).

Thus,

$$\mathfrak{M}_1 = 0.2024$$

The flow through the channel isn’t thermally choked as  $\mathfrak{M}$  doesn’t become 1 by the heat addition.

- (c) After exiting the ducts, the converging-diverging nozzle results in a Mach number given by

$$\frac{A_e}{A^*} = \frac{1}{\mathfrak{M}} \left[ \left( \frac{1 + \frac{\gamma - 1}{2} \mathfrak{M}^2}{\frac{\gamma + 1}{2}} \right)^{\frac{\gamma + 1}{\gamma - 1}} \right]^{1/2}$$

and from the Mach number definition and energy equation (Bernoulli invariant), we can get the exhaust velocity

$$u_e = \mathfrak{M} a = \mathfrak{M} \sqrt{\gamma R T} = \mathfrak{M} \sqrt{\gamma R \frac{T_0}{1 + \frac{\gamma - 1}{2} \mathfrak{M}^2}}$$

**Hydrazine monopropellant thrusters**

Consider in general 2 reactions of the form

$$\sum_j \nu_j A_j = 0$$

and

$$\xi_K A_K + \sum_k \xi_k B_k = 0$$

with the convention that the stoichiometric coefficients “on the left” (reactants) are negative integers.

Clearly,

$$dN_j = \nu_j d\hat{N}$$
$$dN_k = \xi_k d\hat{N}_2$$

where

$$d\hat{N} \equiv \text{how many times reaction 1 takes place}$$
$$d\hat{N}_2 \equiv \text{how many times reaction 2 takes places}$$

since in chemical reactions, chemical species can’t be “broken down” or “destroyed” into smaller parts (we’re not considering nuclear reactions).

Then for the particular chemical species  $A_K$ ,

$$dN_K = \nu_K d\hat{N} \text{ and}$$
$$dN_{K_2} = \xi_{K_2} d\hat{N}_2$$

for chemical species  $K$  to occur in the second reaction. Thus, given the fraction  $X$  that participates in the second chemical reaction, observe, importantly, that

$$X dN_k = -dN_{K_2} = -\xi_{K_2} d\hat{N}_2$$
$$\frac{-X \nu_K}{\xi_{K_2}} d\hat{N} = d\hat{N}_2$$

and so for the other chemical species (molecules) in the second reaction, their “stoichiometric” coefficients change in the following manner:

$$(38) \quad dN_k = \xi_k \frac{X\nu_k}{-\xi_{K_2}} d\hat{N}$$

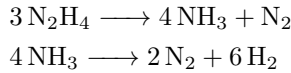
and so writing the first reaction as

$$\sum_j \nu_j A_j = 0 = \sum_{k \neq K} \nu_k A_k + \nu_K A_K$$

then we simply only need to add up reactions 1 and 2, with reaction 2 taking new stoichiometric coefficients according to Eq. [38](#):

$$(39) \quad \boxed{\sum_{k \neq K} \nu_k A_k + \nu_K (1 - X) A_K + \sum_l \frac{\xi_l X \nu_l}{-\xi_{K_2}} B_l = 0}$$

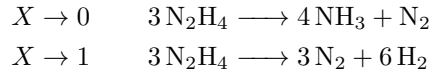
Indeed, consider the (exothermic) decomposition of hydrazine and then the dissociation of ammonia into elevated temperatures:



Then plugging into our formula, Eq. [39](#),

$$\begin{aligned} -3\text{N}_2\text{H}_4 + 4(1 - X)\text{NH}_3 + \text{N}_2 + \frac{2 \cdot X \cdot 4}{4} \text{N}_2 + \frac{6 \cdot X \cdot 4}{4} \text{H}_2 = 0 = -3\text{N}_2\text{H}_4 + 4(1 - X)\text{NH}_3 + (1 + 2X)\text{N}_2 + 6X\text{H}_2 \\ \implies \boxed{\text{N}_2\text{H}_4 \longrightarrow 4(1 - X)\text{NH}_3 + (1 + 2X)\text{N}_2 + 6X\text{H}_2} \end{aligned}$$

Indeed, for the following limits,



#### *adiabatic flame temperature*

Remember that for adiabatic processes,  $Q = 0$ , and so  $Q = 0 = \Delta H = H_{\text{products}} - H_{\text{reactants}}$  and so for adiabatic processes,  $H_{\text{products}} = H_{\text{reactants}}$ , given by

$$(40) \quad \begin{aligned} H_P(\tau_P) &= \sum_{i \in \{\text{products}\}} n_i \left[ \Delta_f \bar{h}_i^0 + \bar{h}_i(\tau_P) - \bar{h}_i^0 \right] \\ H_R(\tau_R) &= \sum_{i \in \{\text{reactants}\}} n_i \left[ \Delta_f \bar{h}_i^0 + \bar{h}_i(\tau_R) - \bar{h}_i^0 \right] \end{aligned}$$

where  $P$  denotes products and  $R$  denotes reactants.

The upper bound of the adiabatic flame temperature is given by no dissociation ( $X = 0$ ). The lower bound of the adiabatic flame temperature is given by the flame temperature with full dissociation ( $X = 1$ ).

EY : 20160121 My question is this: is this because the dissociation of ammonia is endothermic, as the heat of formation of ammonia is  $\Delta H_f^\circ = -45.894 \text{ kJ/mol}$ , and so heat is given off in its formation, and so for full dissociation, the temperature goes down. Is this also because of Le Chatelier’s principle? A clear explanation according to Le Chatelier’s principle would also be appreciated.

Remember that the heat of formation for stable elements is 0.

For the case of  $X = 0$ ,  $Q = 0$  leads to

$$3 \cdot [\Delta_{f;\text{N}_2\text{H}_4}^\circ + 0] = 4 \cdot [\Delta_{f;\text{NH}_3}^\circ + (H - H(T^\circ))_{\text{NH}_3}] + 1 [0 + (H - H(T^\circ))_{\text{N}_2}]$$

For the case of  $X = 1$ ,  $Q = 0$  leads to

$$3 \cdot [\Delta_{f;\text{N}_2\text{H}_4}^\circ + 0] = 3 [0 + (H - H(T^\circ))_{\text{N}_2}] + 6 [0 + (H - H(T^\circ))_{\text{H}_2}]$$

We were given a quadratic polynomial fit to data taken from JANAF (or JANNAF?) tables for  $\text{N}_2, \text{H}_2, \text{NH}_3$ :

$$\begin{array}{ll} \text{N}_2 & H - H(T^\circ) = -8.553 + 27.77\theta + 2.317\theta^2 \\ \text{H}_2 & H - H(T^\circ) = -8.292 + 27.39\theta + 1.586\theta^2 \\ \text{NH}_3 & H - H(T^\circ) = -10.37 + 31.32\theta + 11.77\theta^2 \end{array}$$

where  $\theta = T/1000$  and  $T$  is in Kelvins ( $K$ ).

Solving for a quadratic function is easy with `scipy` using the module `optimize` and the method `newton`. The code is in `combustion.py`.

The upper bound of the adiabatic flame temperature  $T_{ad}$  was found to be

$$T_{ad} = 1644 \text{ K} = 1371 \text{ C} = 2498 \text{ F}$$

and the lower bound of the adiabatic flame temperature was found to be

$$T_{ad} = 866 \text{ K} = 593 \text{ C} = 1099 \text{ F}$$

#### **Gibbs function for a mixture of gases**

Let  $\Sigma$  be a 2-dimensional (i.e.  $\dim \Sigma = 2$ ) manifold representing all possible thermodynamic states of the system (in this case, of an ideal gas). Let  $\Sigma$  have coordinates  $\tau, p$ , representing the temperature  $\tau$ , and pressure  $p$ , of the system, respectively, which could’ve been obtained from (successive) Legendre transformations of (global) coordinates  $U, V$ , where  $U$  is energy.

Let  $Q$ , representing heat applied onto, or heat into, the system, be a 1-form of  $\Sigma$ , i.e.  $Q \in \Omega^1(\Sigma)$ .

Then  $Q$  can be written, in general, as

$$Q = C_p d\tau + \frac{\partial Q}{\partial p} dp$$

Let’s determine  $\frac{\partial Q}{\partial p}$ . From energy conservation (i.e. “First law” of thermodynamics),  $Q = dU + pdV$  (i.e.  $Q$  in (global) coordinates of  $U, V$  for  $\sigma$ ). Next, consider a curve  $\gamma : \mathbb{R} \rightarrow \Sigma$  in  $\Sigma$ , representing a thermodynamic process, such that  $d\tau(\dot{\gamma}) = 0$  which says that the process occurs at constant temperature  $\tau$ . Note that  $\dot{\gamma}$  is a vector field in  $\Sigma$ , i.e.  $\dot{\gamma} \in \mathfrak{X}(\Sigma)$ . Then

$$Q(\dot{\gamma}) = C_p d\tau(\dot{\gamma}) + \frac{\partial Q}{\partial p} dp(\dot{\gamma}) = 0 + \frac{\partial Q}{\partial p} dp(\dot{\gamma}) = dU(\dot{\gamma}) + pdV(\dot{\gamma})$$

For an ideal gas,  $U = \frac{1}{\gamma-1} N\tau$  and so  $dU = \frac{1}{\gamma-1} N d\tau$  and so  $dU(\dot{\gamma}) = \frac{1}{\gamma-1} N d\tau(\dot{\gamma}) = 0$ . It also sufficed to say that for an ideal gas,  $U = U(\tau)$ , i.e.  $U$  is only a function of temperature, and so  $U$  is constant during a process where temperature remains constant.

Then, using  $pV = N\tau$  which must always hold for an ideal gas in equilibrium,

$$V = \frac{N\tau}{p} \implies dV = \frac{N\tau}{-p^2} dp$$

and so

$$\frac{\partial Q}{\partial p} dp = 0 + p \left( \frac{N\tau}{-p^2} \right) dp = \frac{-N\tau}{p} dp = -V dp$$

Thus  $\frac{\partial Q}{\partial p} = -V$ . One can (physically) interpret this as showing how heat, applied to this ideal gas system, can cause our system to be able to do (physical) work.

Then equating  $Q = \tau d\sigma$ ,

$$\tau d\sigma = Q = C_p d\tau - V dp \implies d\sigma = \frac{C_p}{\tau} d\tau - \frac{V}{\tau} dp = \frac{C_p d\tau}{\tau} - \frac{N dp}{p}$$

Taking the integral,

$$\sigma - \sigma_0 = \int_{\tau_0}^{\tau} \frac{C_p d\tau}{\tau} - N \ln \frac{p}{p_0}$$

Let’s talk about, for a mixture of ideal gases, partial pressures  $p_i$  for each of the (kinds of) species  $i$ . It makes (physical) sense that the contribution to the total pressure  $p$  from a single species labeled  $i$  would be proportional to the number of particles  $N_i$  of the species  $i$  (imagine (infinitesimally small) colliding billiard balls). Then  $p_i = \frac{N_i}{N}p$  and so

$$N_i \ln \frac{p_i}{p_0} = N_i \ln \frac{N_i p}{N p_0} = N_i \ln \frac{p}{p_0} + N_i \ln \frac{N_i}{N}$$

For a mixture of ideal gases, there is no interaction energy (by definition of *ideal*), and so entropies are *additive* [14]. So

$$\sigma = \sum_i \sigma_i = \sum_i \left( \int_{\tau_0}^{\tau} \frac{C_{pi} d\tau}{\tau} + \sigma_{i0} \right) - N \ln \frac{p}{p_0} - \sum_i N_i \ln \frac{N_i}{N}$$

$\sigma$  in its current state is unitless. Multiply by the Boltzmann constant  $k_B$  and Avogadro’s number  $N_A$  as appropriate to obtain the desired units for entropy  $S$ .

12.4.1. *ig thermo.xls; Dealing with tables.* Doing a search on Google for `ig_thermo.xls` yields this link:

[http://www.google.com/url?q=http://shepherd.caltech.edu/EDL/projects/pde/ig\\_thermo.xls&sa=U&ved=0ahUKEwj2ta8x7nKAhVB8GMKHa-tCBYQFggUMAA&usg=AFQjCNGcu2ryEOaIZqwumHooxDjcxTf3gQ](http://www.google.com/url?q=http://shepherd.caltech.edu/EDL/projects/pde/ig_thermo.xls&sa=U&ved=0ahUKEwj2ta8x7nKAhVB8GMKHa-tCBYQFggUMAA&usg=AFQjCNGcu2ryEOaIZqwumHooxDjcxTf3gQ)

However, ammonia is not to be found in `ig_thermo.xls` *off the Web*.

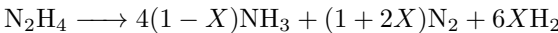
Let’s consider the work-flow or process of scraping the NIST website for JANAF table data. Starting from the **NIST-JANAF Thermochemical Tables** webpage, one can enter, as a string, the CAS number, chemical formula, or compound name, or search by the periodic table. However, using the periodic table, for element N or H, ammonia wasn’t to be found (as of 20160120). However, entering NH<sub>3</sub> or Ammonia (the search text box “field” is case-sensitive), leads to the ammonia gas JANAF table. Then one clicks on the link for “**view**” for the HTML-formatted table and then follow the bottom link for the tab-separated table.

So it would be nice to do the following:

- Obtain a *comprehensive* list of all available compounds of this NIST JANAF database and with that
- scrape all the tab-separated tables.

A comprehensive list appears to be in this pdf link for **Volume 1** of the JANAF tables. Perhaps this pdf could be “scraped” to obtain a list of compounds in text format?

12.4.2. *entropy of a mixture of ammonia, nitrogen, and hydrogen.* Recall that



and that we had derived

$$S(T, P) = \sum_i N_i \left( \int_{T_0}^T c_{pi} \frac{dT}{T} + s_{i0} \right) - NR \ln \frac{P}{P_0} - R \sum_i N_i \ln \frac{N_i}{N}$$

The values we want are for  $P = 1$  atm and  $P = 10$  atm

Note that

$$N = 4(1 - X) + (1 + 2X) + 6X = 4X + 5$$

and

$$\sum_i N_i \ln \frac{N_i}{N} = 4(1 - X) \ln \left( \frac{4(1 - X)}{4X + 5} \right) + (1 + 2X) \ln \left( \frac{1 + 2X}{4X + 5} \right) + 6X \ln \left( \frac{6X}{4X + 5} \right)$$

### 13. EQUILIBRIUM FLOW VS. FROZEN FLOW FOR NOZZLE FLOW AND THE EXAMPLE OF THE SPACE SHUTTLE MAIN ENGINE (SSME)

- Remember that from Newton’s 2nd law, the so-called momentum equation, is

$$\rho \frac{Du}{Dt} = \text{div}(\mathbf{T}) \text{ or } \rho \left[ \frac{\partial \mathbf{u}}{\partial t} + u^j \frac{\partial \mathbf{u}}{\partial x^j} \right] = \text{div}(\mathbf{T})$$

Let

$$k := \frac{1}{2} \rho u^2$$

which is a scalar quantity that depends on time  $t$  and spatial coordinates  $\mathbf{x} = (x^1 \dots x^n) \in N$  and represents the density of kinetic energy of the bulk fluid at time  $t$  and spatial point  $\mathbf{x}$ .

Clearly, product rule still applies when using the material derivative  $\frac{D}{Dt} := \frac{\partial}{\partial t} + u^j \frac{\partial}{\partial x^j}$ :

$$\frac{D}{Dt} k = \frac{D}{Dt} \frac{1}{2} \rho u^2 = \frac{1}{2} \left[ u^2 \frac{D\rho}{Dt} + \rho \left( \frac{D}{Dt} u^2 \right) \right] = \frac{1}{2} \left[ u^2 (-\rho \text{div} u) + \rho \left( \frac{D}{Dt} u^2 \right) \right]$$

If, say,  $g = 1$ , then  $u^2 = u_i u^i = (u^i)^2$ , and so

$$\rho \frac{D}{Dt} u^2 = 2 \rho u^i \frac{D}{Dt} u^i = 2 u^i \text{div}(T^i)$$

and so in this case,

$$\frac{D}{Dt} k = \frac{-1}{2} \rho u^2 \text{div}(u) + u^i \text{div}(T^i) = -k \text{div}(u) + u^i \text{div}(T^i)$$

and since

$$u^j \frac{\partial k}{\partial x^j} + k \text{div}(u) = \text{div}(ku)$$

then we could also write

$$\frac{D}{Dt} k + k \text{div}(u) = \frac{\partial k}{\partial t} + \text{div}(ku) = u^i \text{div}(T^i)$$

What about the case where  $g \neq 1$  but could be dependent upon time and space?

Invoke the “Fundamental Theorem on Time-Dependent Flows” (Theorem 9.48 in Chapter 9 Integral Curves and Flows, pp. 237 of Lee (2012) [16], which guarantees the existence of a time-dependent flow  $\phi$  for our velocity vector field  $\mathbf{u}$ .

So for our time-dependent velocity vector field,  $\mathbf{u} \in \mathfrak{X}(\mathbb{R} \times N)$ , i.e.

$$\mathbf{u} : \mathbb{R} \times N \rightarrow TN$$

$$\mathbf{u} = \mathbf{u}(t, \mathbf{x}) = u^i(t, \mathbf{x}) \frac{\partial}{\partial x^i}$$

there exists a flow, given initial conditions,  $\phi = \phi(t) \in N$ , such that

$$\phi(t) = (x^1(t) \dots x^n(t))$$

$$\dot{\phi}(t) = (\dot{x}^1(t) \dots \dot{x}^n(t)) = (u^1 \dots u^n)$$

Thus, again emphasizing on any Riemannian manifold  $(M, g)$ , equipped with metric  $g$ ,

$$u^2 = g_{ij} u^i u^j = g_{ij} \dot{\phi}^i \dot{\phi}^j$$

Consider any time and space dependent scalar quantitiy,  $f$ ,  $f \in C^\infty(\mathbb{R} \times N)$ , i.e.  $f : \mathbb{R} \times N \rightarrow \mathbb{R}$ . On a flow  $\phi$  for a time-dependent vector field  $\mathbf{u}$  generating this flow,

$$\begin{aligned} \frac{d}{dt} f &= \frac{d}{dt} f(t, \mathbf{x}) = \frac{d}{dt} f(t, \mathbf{x}(t)) = \frac{\partial}{\partial t} f + \frac{\partial f}{\partial x^i} \dot{x}^i = \frac{\partial}{\partial t} f + u^i \frac{\partial f}{\partial x^i} = \\ &= \frac{D}{Dt} f \end{aligned}$$

Thus, in this scenario, then  $\forall f \in C^\infty(\mathbb{R} \times N)$ , the material derivative  $\frac{D}{Dt}$  of  $f$  is equal to the total time derivative  $\frac{d}{dt}$ .

Then for scalar quantity  $u^2 = u^2(t, \mathbf{x}) = g_{ij}(\phi(t)) \dot{\phi}^i \dot{\phi}^j$ ,

$$\frac{d}{dt} u^2 = \frac{\partial g_{ij}}{\partial x^k} \dot{\phi}^k \dot{\phi}^i \dot{\phi}^j + 2g_{ij} \dot{\phi}^i \ddot{\phi}^j = 2g_{ks} \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \dot{\phi}^s + 2g_{sk} \dot{\phi}^s \dot{\phi}^k = 2g_{ks} \dot{\phi}^s \left[ \ddot{\phi}^k + \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \right]$$

where the identity

$$\frac{\partial g_{ij}}{\partial x^k} \dot{\phi}^k \dot{\phi}^i \dot{\phi}^j = 2g_{ks} \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \dot{\phi}^s$$

was used, which is in the proof for Theorem 3.25 on pp. 49 of section 3.8 The natural Lagrangian on manifolds of Calin and Chang (2005) [17]. I want to remark that Calin and Chang (2005) [17] has explicit computations showing this identity to be the case locally, and at least for me, I don't see many other sources that are this thorough with presenting explicit computations.

Noting that

$$\nabla_u u = \nabla_{u^j \frac{\partial}{\partial x^j}} u = u^j \nabla_{\frac{\partial}{\partial x^j}} u = u^j \left( \frac{\partial u^k}{\partial x^j} + \Gamma_{ij}^k u^i \right) \frac{\partial}{\partial x^k} = \left( \ddot{\phi}^k + \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \right) \frac{\partial}{\partial x^k}$$

and so

$$g_{ks} \dot{\phi}^s \left[ \ddot{\phi}^k + \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \right] = \langle u, \nabla_u u \rangle \implies \frac{d}{dt} u^2 = 2 \langle u, \nabla_u u \rangle$$

Thus, in general (on a Riemannian manifold),

$$(42) \quad \boxed{\frac{\partial k}{\partial t} + \text{div}(ku) = \frac{Dk}{Dt} + k \text{div}(u) = \rho \langle u, \nabla_u u \rangle = \langle u, \text{div}(T) \rangle}$$

Now when following a volume element  $\text{vol}^n$ , along the flow of the fluid, one should be concerned that in a “lab frame”, the volume itself can change with time:

$$\begin{aligned} \dot{V} &= \frac{d}{dt} \int_V \text{vol}^n = \int_V \mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}} \text{vol}^n = \int_V 0 + \mathcal{L}_{\mathbf{u}} \text{vol}^n = \int_V (\mathbf{d}i_{\mathbf{u}} + i_{\mathbf{u}} \mathbf{d}) \text{vol}^n = \int_V \mathbf{d}i_{\mathbf{u}} \text{vol}^n + 0 = \\ &= \int_{\partial V} i_{\mathbf{u}} \text{vol}^n = \int_{\partial V} u^j dS_j \end{aligned}$$

Indeed, if the time-dependent velocity vector field  $\mathbf{u}$  changes at the boundary  $\partial V$ , then so does our volume, in this frame.

The correct way out of this conundrum is to switch to the frame where the bulk velocity of the fluid is 0, so it's at rest in this frame. This is paramount to choosing  $\mathbf{u} = 0$  above, and so  $\dot{V} = 0$ .

The strategy is this: one should recognize that number of particles  $N$ , or in general, number of particles for species  $i$ ,  $N_i$  and entropy  $\sigma$  are invariants under Galilean (or Lorentz) transformation: observers, whether in the lab frame, or the frame in which the bulk velocity of the fluid is 0 so the bulk fluid flow is at rest, or any other frame, must agree upon these numbers (that can be objectively measured). So we've fixed our “control volume” that remains constant in time in this fluid-at-rest frame, and count the number of particles that come in, and the entropy.

For only 1 kind of particle, let

$$\check{\sigma} \equiv \sigma/N \text{ entropy per particle}$$

$$\check{\epsilon} \equiv E/N \text{ energy per particle}$$

Thus

$$\tau d\sigma = \tau d(N\check{\sigma}) = \tau N d\check{\sigma} + \tau \check{\sigma} dN$$

where

$$\tau N d\check{\sigma} \text{ is entropy change due to change in entropy per particle, i.e. } \quad \mathbf{conduction term}$$

$$\tau \check{\sigma} dN \text{ is entropy change due to change in number of particles, i.e. } \quad \mathbf{convection term}$$

Generalizing this to  $\mathcal{N}$  species, indexed by  $i = 1 \dots \mathcal{N}$ ,

$$\check{\sigma}_i \equiv \sigma_i/N_i \text{ entropy per particle of species } i$$

Assuming that entropies are additive (which is a reasonable assumption if the interaction between species is negligible [14]), then for total entropy of the system in volume  $V$  in this fluid-at-rest frame (which will be denoted with prime as needed),

$$\sigma = \sum_i \sigma_i = \sum_i \check{\sigma}_i N_i$$

So

$$\tau d\sigma = \tau \sum_i N_i d\check{\sigma}_i + \tau \sum_i \check{\sigma}_i dN_i$$

where

$$\tau \sum_i N_i d\check{\sigma}_i \text{ is entropy change due to change in entropy per particle for each species, i.e. } \quad \mathbf{conduction term}$$

$$\tau \sum_i \check{\sigma}_i dN_i \text{ is entropy change due to change in number of particles for each species, i.e. } \quad \mathbf{convection term}$$

Using the above result for the decomposition of  $\tau d\sigma$  into a conduction term and convection term, then

$$(43) \quad \begin{aligned} dE &= \tau d\sigma - pdV + \sum_i \mu_i dN_i = \tau \sum_i N_i d\check{\sigma}_i + \tau \sum_i \check{\sigma}_i dN_i - pdV + \sum_i \mu_i dN_i = \\ &= \tau \sum_i N_i d\check{\sigma}_i + \sum_i (\tau \check{\sigma}_i + \mu_i) dN_i - pdV = \tau \sum_i N_i d\check{\sigma}_i + \sum_i \check{h}_i dN_i - pdV = \tau \sum_i N_i d\check{\sigma}_i + \sum_i \check{h}_i dN_i \end{aligned}$$

since  $dV = 0$  for any thermodynamic process in this constant volume and that

$$G_i = F_i + p_i V = E_i - \tau \sigma_i + p_i V = H_i - \tau \sigma_i = \mu_i N_i \implies \check{h}_i = \mu_i + \tau \check{\sigma}_i$$

since Gibbs free energy  $G_i$  for each species  $i$  should be additive, partial pressures  $p_i$  for each species  $i$  should be additive to the total pressure  $p$ , as forces are additive, and from the definition of enthalpy  $H_i = E_i + p_i V$  for each species  $i$ .

For the case of 1 species, and for this volume at rest in the fluid-at-rest frame,  $V$ , then defining

$$\begin{aligned} h &:= \frac{H}{V} \\ s &:= \frac{\sigma}{V} \\ \epsilon &:= \frac{E}{V} \end{aligned}$$

and so from Eq. 43,

$$d\epsilon = \tau \frac{N}{V} d\left(\frac{\sigma}{N}\right) + \frac{H}{N} d\left(\frac{N}{V}\right) = \tau \frac{N}{V} \left( \frac{d\sigma}{N} + -\frac{\sigma dN}{N^2} \right) + \frac{H}{N} \frac{d\sigma}{M} = \tau ds - \frac{\tau s}{\rho} d\rho + \frac{h d\rho}{\rho} = \tau ds + \frac{h - \tau s}{\rho} d\rho$$

We should denote with prime symbols (') the quantities that transform under change of observer's frame, and is specific, in value, to this fluid-at-rest frame:

$$d\epsilon' = \tau ds + \frac{h' - \tau s}{\rho} d\rho$$

For many species, then generalizing the above,

$$d\epsilon' = \tau \sum_i \frac{N_i}{V} d\left(\frac{\sigma_i}{N_i}\right) + \sum_i \frac{H_i}{N_i} d\frac{N_i}{V} = \tau \sum_i \frac{N_i}{V} \left( \frac{d\sigma_i}{N_i} + -\frac{\sigma_i dN_i}{N_i^2} \right) + \sum_i \frac{h'_i}{\rho_i} d\rho_i = \tau ds + \sum_i \frac{h'_i - \tau s_i}{\rho_i} d\rho_i$$

Thus, in this fluid-at-rest frame,

$$(44) \quad \boxed{d\epsilon' = \tau ds + \sum_i \frac{h'_i - \tau s_i}{\rho_i} d\rho_i}$$

Next, use the fact that for a fluid in the “lab frame”, the current density for the total energy  $E$  is carried by  $(k + h')\mathbf{u}$  and *not*  $(k + \epsilon')\mathbf{u}$ . This comes from the fact that for compressible flows, the Bernoulli invariant is  $k + h' = \frac{1}{2}\rho u^2 + h'$ , and not  $k + \epsilon'$ . The physical interpretation is that the enthalpy  $h'$  is needed to account for all convection terms, to account for the “energy balancing” needed, done by the compression of the volume by  $vdN_i := V/N_i dN_i$  to return to the initial volume  $V$  in the fluid-at-rest frame, when particles of species  $i$  flow through the volume  $V$ , convection.

Thus

$$(45) \quad \frac{\partial \epsilon}{\partial t} + \text{div}((k + h')\mathbf{u}) = 0 \text{ or } \frac{\partial \epsilon}{\partial t} + \text{div}(ku) = -\text{div}(h'u)$$

With all these ingredients, begin with the account of all energies in the lab frame, for a fluid:

$$E = KE + U \implies \epsilon = k + \epsilon' \implies \frac{\partial \epsilon}{\partial t} = \frac{\partial k}{\partial t} + \frac{\partial \epsilon'}{\partial t}$$

From Eq. 42, to substitute into  $\frac{\partial k}{\partial t}$ ,

$$\frac{\partial \epsilon}{\partial t} = -\text{div}(ku) + \langle u, \text{div}T \rangle + \frac{\partial \epsilon'}{\partial t}$$

From Eq. 44, we can consider a thermodynamic process where we vary the time, generated by vector  $\frac{\partial}{\partial t}$ , so that we have

$$\frac{\partial \epsilon'}{\partial t} = \tau \frac{\partial s}{\partial t} + \sum_i \frac{h'_i - \tau s_i}{\rho_i} \frac{\partial \rho_i}{\partial t}$$

and so

$$\frac{\partial \epsilon}{\partial t} = -\text{div}(ku) + \langle u, \text{div}T \rangle + \tau \frac{\partial s}{\partial t} + \sum_i \frac{h'_i - \tau s_i}{\rho_i} \frac{\partial \rho_i}{\partial t} = -\text{div}(ku) - \text{div}(h'u)$$

where in the last equality, we used Eq. 45. Eliminating  $-\text{div}(ku)$  from both sides, noting that

$$\text{div}(h'u) = \frac{1}{\sqrt{g}} \frac{\partial(h'u^j \sqrt{g})}{\partial x^j} = u^j \frac{\partial h'}{\partial x^j} + h' \text{div}(u)$$

and using  $T = -p1$ , the case where the stress tensor is only the isotropic (same in all directions) pressure  $p$ , then

$$(46) \quad -u^j \frac{\partial p}{\partial x^j} + \tau \frac{\partial s}{\partial t} + \sum_i \frac{h'_i - \tau s_i}{\rho_i} \frac{\partial \rho_i}{\partial t} = -u^j \frac{\partial h'}{\partial x^j} - h' \text{div}(u)$$

Now we'll need to use a mass conservation law for each of the species  $i$ , which I'm not sure is valid (please correct me if I'm wrong):

$$\frac{\partial \rho_i}{\partial t} + \text{div}(\rho_i u) = 0 \implies \frac{\partial \rho_i}{\partial t} = -\rho_i \text{div}u - u^j \frac{\partial \rho_i}{\partial x^j}$$

and so

$$\sum_i \frac{h'_i - \tau s_i}{\rho_i} \frac{\partial \rho_i}{\partial t} = \sum_i -h'_i \text{div}u + \tau s_i \text{div}u - \frac{h'_i - \tau s_i}{\rho_i} u^j \frac{\partial \rho_i}{\partial x^j} = -h' \text{div}u + \tau s \text{div}u - \sum_i \frac{h'_i - \tau s_i}{\rho_i} u^j \frac{\partial \rho_i}{\partial x^j}$$

Plugging this into Eq. 46,

$$-u^j \frac{\partial p}{\partial x^j} + \tau \frac{\partial s}{\partial t} + -h' \text{div}u + \tau s \text{div}u - \sum_i \frac{h'_i - \tau s_i}{\rho_i} u^j \frac{\partial \rho_i}{\partial x^j} = -u^j \frac{\partial h'}{\partial x^j} - h' \text{div}(u)$$

Noting that by definition,  $h' := H'/V = \epsilon' + p$  Using Eq. 44 one more time so that

$$\frac{\partial \epsilon'}{\partial x^j} u^j = \left[ \tau \frac{\partial s}{\partial x^j} + \sum_i \frac{h'_i - \tau s_i}{\rho_i} \frac{\partial \rho_i}{\partial x^j} \right] u^j$$

and thus

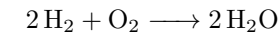
$$\tau \frac{\partial s}{\partial t} + \tau s \text{div}u + \tau u^j \frac{\partial s}{\partial x^j} = 0 \implies \boxed{\frac{\partial s}{\partial t} + \text{div}(su)}$$

Thus, entropy is conserved over time. Indeed

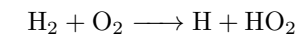
$$\dot{S} = \frac{d}{dt} \int_V s \text{vol}^n = \int_V \left( \frac{\partial s}{\partial t} + \text{div}(su) \right) \text{vol}^n$$

Also notice that I did not assume incompressibility. This, and the statements above, are true for compressible and incompressible flow. Also, note that all the above was proven for any Riemannian manifold  $N$  representing space.

(2) I supposed the reaction was



But the reaction listed for `h2o2.cti` or `h2o2_highT.cti` is



#### 14. LIQUID-VAPOR EQUILIBRIUM

From [wikipedia](#)'s article on "Clausius-Clapeyron relation", start from the definition of particle diffusion equilibrium, that

$$\begin{aligned} \mu_v &= \mu_l \\ d\mu_v &= d\mu_l \end{aligned}$$

Using the Gibbs-Duhem relation for each side of the above, that  $Nd\mu = Vdp - \sigma d\tau$ , then

$$\begin{aligned} d\mu_v &= v_v dp - s_v d\tau \\ d\mu_l &= v_l dp - s_l d\tau \end{aligned}$$

Thus

$$(v_v - v_l)dp = (s_v - s_l)d\tau \implies \boxed{\frac{dp}{d\tau} = \frac{s_v - s_l}{v_v - v_l}}$$

Using the definition of enthalpy  $H := U + pV$ , so that

$$dH = dU + Vdp + pdV = \tau d\sigma + Vdp \text{ or } \frac{1}{\tau}(dH - Vdp) = d\sigma$$

where  $\tau d\sigma = dU + pdV$  was used.

Considering the thermodynamic process at constant pressure  $dp = 0$  that begins with a single molecule being in a liquid state and changing into vapor state, then

$$\frac{1}{\tau}(\Delta H_v) = \sigma_v - \sigma_l$$

Thus

$$\frac{dp}{dT} = \frac{\Delta H_v}{T(V_v - V_l)} \approx \frac{\Delta H_v}{TV_v}$$

since it's reasonable to assume that the volume of the vapor molecule is much greater than when the molecule is in a liquid state.

Using the ideal gas law for a single molecule,  $pV = \tau$ , then

$$\begin{aligned} \frac{dp}{dT} &= \frac{p\Delta H_v}{T^2} \text{ or } \frac{dp}{p} = \frac{\Delta H_v dT}{T^2} \\ \implies p(T) &= p_0 \exp\left(\frac{-\Delta H_v}{T}\right) \end{aligned}$$

#### Part 5. Basic Feeling

##### 15. BOX WITH A HOLE ROCKET; BOTTLED (BOX) ROCKET

Recall Ch. 14: Kinetic Theory, Section "Kinetic Theory of the Ideal Gas Law" of Kittel and Kroemer [8].



**Kinetic Theory of the Ideal Gas Law.** Consider molecule strike unit area of wall.

Let  $v_z \equiv$  velocity component normal to plane of wall.

Suppose molecules, of mass  $M$ , reflected specularly (mirror-like) from wall,

$$\Delta p_z = -2M|v_z|$$

Let  $a(v_z)dv_z$ , number of molecules per unit volume with  $z$ -component of velocity between  $v_z$  and  $v_z + dv_z$ .

$$\int a(v_z)dv_z = \frac{N}{V} = n$$

$a(v_z)v_zdv_z$  number of molecules in  $(v_z, v_z + dv_z)$  velocity range that strike unit area of wall in (per) unit time

$$\text{pressure } p = \int_0^\infty 2Mv_z a(v_z)v_z dv_z = M \int_{-\infty}^\infty v_z^2 a(v_z)dv_z = Mn\langle v_z^2 \rangle$$

$\frac{1}{2}M\langle v_z^2 \rangle = \frac{1}{2}\tau$  by equipartition of energy (Ch.3)

$$p = nM\langle v_z^2 \rangle = n\tau = \frac{N\tau}{V}$$

*Maxwell Distribution of Velocities.* cf. Ch.6. distribution function of ideal gas  $f(\epsilon_n) = \lambda \exp\left(\frac{-\epsilon_n}{\tau}\right)$

Recall, Ch. 6, Sec. “Classical Limit” of Kittel and Kroemer [8], **an ideal gas is defined as a system of free noninter-acting particles in the classical regime.**

$f(\epsilon) \equiv$  average occupancy of an orbital at energy  $\epsilon$

$\epsilon \equiv$  energy of orbital occupied by 1 particle; not energy of system of  $N$  particles

Fermi-Dirac and Bose-Einstein distribution  $f(\epsilon) = \frac{1}{\exp[(\epsilon - \mu)/\tau] \pm 1}$

In order for  $f(\epsilon) \ll 1 \forall$  orbitals,  $\exp[(\epsilon - \mu)/\tau] \gg 1 \forall \epsilon$ .

$$\implies f(\epsilon) \simeq \exp[(\mu - \epsilon)/\tau] = \lambda \exp(-\epsilon/\tau) \quad \lambda \equiv \exp\left(\frac{\mu}{\tau}\right)$$

$f(\epsilon)$ , average occupancy of orbital of energy  $\epsilon$ , is classical distribution function.

particle in a box:  $\epsilon_n = \frac{1}{2M} \left(\frac{\pi n}{L}\right)^2$  (for, recall  $\frac{1}{2M} \left(\frac{1}{i}\partial\right)^2 \psi = \frac{-1}{2M} \partial^2 \psi = E\psi$ )

number of orbitals in range of quantum number  $(n, n + dn)$ , probability such orbital is occupied

$$\left(\frac{1}{2}\pi n^2 dn\right) f(\epsilon_n) = \frac{1}{2}\pi n^2 \lambda \exp(-\epsilon_n/\tau) dn$$

$$\frac{1}{2}Mv^2 = \frac{1}{2M} \left(\frac{\pi n}{L}\right)^2 \quad \text{or } n^2 = \frac{(ML)^2}{\pi^2} v^2 \quad \text{or } n = \frac{MLv}{\pi}$$

Consider system of  $N$  particles in volume  $V$ .

Let  $NP(v)dv$  number of atoms with velocity magnitude in range  $dv$  at  $v$

$$NP(v)dv = \frac{1}{2}\pi n^2 \lambda \exp(-\epsilon_n/\tau) \frac{dn}{dv} dv = \frac{1}{2}\pi \lambda \left(\frac{ML}{\pi}\right)^3 v^2 \exp\left(\frac{-Mv^2}{2\tau}\right) dv$$

cf. Ch. 6,  $\lambda = \frac{n}{n_Q} = \frac{N}{L^3} \left(\frac{2\pi\hbar^2}{M\tau}\right)^{3/2}$

$$\frac{1}{2}\pi \frac{N}{L^3} \left(\frac{2\pi}{M\tau}\right)^{3/2} \left(\frac{ML}{\pi}\right)^3 = 4\pi N \left(\frac{N}{2\pi\tau}\right)^{3/2}$$

$$(47) \quad \implies P(v) = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} v^2 \exp\left(\frac{-Mv^2}{2\tau}\right)$$

$P(v)$  is **Maxwell velocity distribution**,  $P(v)dv$  is probability particle has speed in  $dv$  at  $v$

*Experimental verification.* velocity distribution of atoms which exit from slit of oven.

exit beam weighted in favor of atoms of high velocity at expense of those at low velocity.

weight factor is velocity component  $v \cos \theta$  normal to plane of hole.

$$\begin{aligned} \int (\cos \theta) drrd\varphi r \sin \theta d\theta &= \left(\frac{1}{3}(2\pi)R^3\right) \int \cos \theta \sin \theta d\theta = \left(\frac{2\pi}{3}R^3\right) \int \frac{\sin 2\theta d\theta}{2} = \frac{2\pi R^3}{3} \left(\frac{-\cos 2\theta}{4}\right) \Big|_0^{\pi/2} = \\ &= \frac{2\pi R^3}{3} \left(\frac{1+1}{4}\right) = \frac{2\pi R^3}{3} \left(\frac{1}{2}\right) \end{aligned}$$

Probability atom leaves hole will have velocity in  $(v, v + dv)$  :  $P_{\text{beam}}(v)dv$

$$P_{\text{beam}}(v) \propto v P_{\text{Maxwell}} \propto v^3 \exp\left(\frac{-Mv^2}{2\tau}\right)$$

with, recall  $P_{\text{Maxwell}} = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} v^2 \exp\left(\frac{-Mv^2}{2\tau}\right)$

$P_{\text{beam}}$  distribution of transmission through a hole is Maxwell transmission distribution.

$$(48) \quad \begin{aligned} \langle v_{\text{out}} \rangle &= \int \int_0^{\pi/2} v \cos \theta \sin \theta d\theta P_{\text{Maxwell}}(v) dv = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \frac{1}{2} \int_0^\infty v^3 \exp\left(\frac{-M}{2\tau}v^2\right) dv = \\ &= 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \frac{1}{2} \left(\frac{1}{-2\alpha}\right) \left[0 - \frac{1}{\alpha}\right] = 4\pi \left(\frac{M}{2\pi\tau}\right)^{3/2} \frac{1}{4} \left(\frac{4\tau^2}{M^2}\right) = \boxed{\frac{2^{1/2}}{\pi^{1/2}} \left(\frac{\tau}{M}\right)^{1/2}} \end{aligned}$$

for (doing the integration by hand)

$$(e^{-\alpha v^2})' = -2\alpha v e^{-\alpha v^2}$$

$$(v^2 e^{-\alpha v^2})' = -2\alpha v^3 e^{-\alpha v^2} + 2v e^{-\alpha v^2}$$

$$(v^2 e^{-\alpha v^2} + \frac{e^{-\alpha v^2}}{\alpha})' = -2\alpha v^3 e^{-\alpha v^2}$$

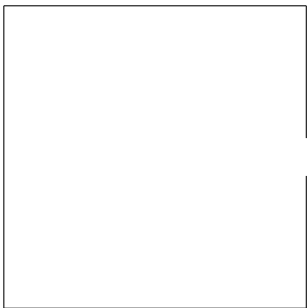
Armed with this mean velocity  $\langle v_{\text{out}} \rangle$  out of a hole of a box, we want the *thrust* that results on a box if we had a box of air, at some pressure, and at some temperature, and then we punch a hole at one end.

What’s happening? Air that’s swirling around the box is then accelerated out of the hole. There’s fluid flow out. For low enough velocities, use *Bernoulli’s equation*, and assume at the starting point of the air’s streamline, the velocity is 0:

$$\begin{aligned} \frac{1}{2}u^2 + \frac{p}{\rho} &= \frac{1}{2}u_f^2 + \frac{p_f}{\rho_f} \\ \implies \frac{p}{\rho} &= \frac{1}{2}u_f^2 + \frac{p_f}{\rho_f} \end{aligned}$$

$\rho$  is really  $\frac{MN}{V}$  and by the ideal gas law (still applies),  $pV = N\tau$  and  $\frac{p}{\tau} = \frac{N}{V}$  but the point is the gas density didn’t change much.





Thus,

$$\frac{p - p_f}{\rho} = \frac{1}{2} u_f^2$$

The thrust is going to be given by the difference in pressure against the walls, the wall in front of the box opposite to the wall with a hole. *You don't need the area of the hole.* This thrust is  $(p - p_f)$

$$(p - p_f)A = \rho \frac{1}{2} u_f^2 A = \frac{MN}{V} \frac{1}{2} u_f^2 A = \frac{MN}{L} \frac{1}{2} u_f^2$$

Now from Eq. 48,  $u_f^2 = \frac{2}{\pi} \frac{\tau}{M}$ , and so

$$(p - p_f)A = \frac{N\tau}{\pi L} = \frac{pV}{\pi L} = \frac{pL^2}{\pi}$$

Thus, the thrust on this box is given by

$$F_{\text{thrust}} = \frac{N\tau}{\pi L} = \frac{pL^2}{\pi}$$

## Part 6. Combustion

### 16. MASS FRACTION VS. MOLE FRACTION VS. MOLECULAR MASS I.E. “MOLECULAR WEIGHT”

This follows Powers (2014) [12], Section 2.1. Some general issues, Chapter 2. Gas Mixtures, as some of the following concepts need to be clarified.

Consider a mixture of  $\mathcal{N} \in \mathbb{Z}^+$ , each “a pure substance”.

Let the total mass of the mixture be  $M$  such that ( $\equiv$  s.t.)

$$M = \sum_{j=1}^{\mathcal{N}} m_j$$

where  $m_j \equiv$  is the mass of the  $j$ th substance.

Let  $N$  be the total number of *particles* s.t.

$$N = \sum_{j=1}^{\mathcal{N}} N_j$$

where  $N_j$  is the (total) number of particles of the  $j$ th substance.

Define

$$N_j/N_A =: n_j$$

with  $n_j$  be the number of moles of the  $j$ th substance, and where  $N_A$  is *Avogadro's number*, which really is a conversion *factor* between units of moles to units for number of particles.

Then

$$n = \sum_{j=1}^{\mathcal{N}} n_j$$

is the (total) number of moles of the mixture, or the “total number” (pp. 73, Powers (2014)) [12]. Define

$$(49) \quad Y_j := \frac{m_j}{M} \quad (\text{mass fraction})$$

$Y_j$  to be the mass fraction.

Define

$$(50) \quad X_j := \frac{n_j}{n} = \frac{N_j}{N} \quad (\text{mole fraction or particle fraction})$$

$X_j$  to be the mole fraction or particle fraction or just fraction.

Define the *molecular mass* or “molecular weight”  $\hat{m}_j$  to be

$$(51) \quad \hat{m}_j := \frac{m_j}{n_j} = \frac{m_j}{N_j} N_A$$

Then the mass fraction  $Y_j$  in terms of molecular mass and in terms of mole fraction, respectively, is

$$Y_j = \frac{m_j}{m} = \frac{\hat{m}_j n_j}{\sum_{\xi=1}^{\mathcal{N}} n_{\xi} \hat{m}_{\xi}} = \frac{\hat{m}_j X_j}{\sum_{\xi=1}^{\mathcal{N}} X_{\xi} \hat{m}_{\xi}}$$

Mole fraction in terms of mass fraction is

$$X_j = \frac{n_j}{n} = \frac{m_j / \hat{m}_j}{\sum_{k=1}^{\mathcal{N}} m_k / \hat{m}_k} = \frac{Y_j / \hat{m}_j}{\sum_{k=1}^{\mathcal{N}} Y_k / \hat{m}_k}$$

The mean molecular mass (of the mixture) is

$$\widehat{M} := \frac{M}{n} = \sum_{j=1}^{\mathcal{N}} \frac{n_j}{n} \hat{m}_j = \sum_{j=1}^{\mathcal{N}} X_j \hat{m}_j$$

Example 2.1 of Powers (2014) [12]

I will now follow the second edition (2000) of Turns [13] (note that there is a third edition for 2011, but I only have the second edition available to me).

On the end of subsection “Ideal-Gas Mixtures” (pp. 15) of Chapter 2 Combustion and Thermochemistry of Turns (2000) [13], recall the Sackur-Tetrode equation for monatomic ideal gas [8]

$$\sigma = N \left[ \ln \left( \frac{n_Q}{n} \right) + \frac{d}{2} + 1 \right] \quad \text{where} \quad n_Q = \left( \frac{M\tau}{2\pi\hbar^2} \right)^{d/2}$$

Now  $\forall j$ th substance, let  $n_{Q;j} := \left( \frac{M_j\tau}{2\pi\hbar^2} \right)^{d/2}$ , since each substance will have a different mass  $M_j$  for each of its molecules, and

let  $n = n_j = \frac{N_j}{V} = \frac{p_j}{\tau}$ . Then

$$\sigma_j = N_j \left[ \ln(n_{Q;j}\tau) - \ln p_j + \frac{d}{2} + 1 \right]$$

where  $\mathcal{N}$  is the total number of substances in the mixture.

If we were to add these entropies up for each of the substances, assuming the additivity of entropies, and thus assuming that the interaction between the disparate substances is negligible,

$$\sum_{j=1}^{\mathcal{N}} \sigma_j = \sum_{j=1}^{\mathcal{N}} N_j [\ln n_{Q;j} - \ln p_j] + N_j [\ln \tau + \frac{d}{2} + 1] = \sum_{j=1}^{\mathcal{N}} N_j (\ln n_{Q;j} - \ln p_j) + N (\ln \tau + 1 + \frac{d}{2})$$

If we wanted to also include internal degrees of freedom (e.g. vibrational and rotational degrees of freedom), the expression for the entropy becomes (derived in [thermo.pdf](#), notes on thermodynamics by me)

$$\sigma = N \left[ \ln \left( \frac{n_Q}{n} \right) + \left( \frac{d}{2} + 1 \right) \right] + \sigma_{\text{int}}$$

and so

$$\sum_{j=1}^{\mathcal{N}} \sigma_j = \sum_{j=1}^{\mathcal{N}} N_j \ln n_{Q;j} + \sigma_{\text{int};j} - N_j \ln p_j + N (\ln \tau + 1 + \frac{d}{2})$$

with

$$\sigma_j = N_j \ln n_{Q;j} + \sigma_{\text{int};j} - N_j \ln p_j + N_j (\ln \tau + 1 + \frac{d}{2})$$

We can make a *formal* definition, with respect to a reference pressure  $p_{\text{ref}}$ , of  $\sigma_j$  so that

(52) 
$$\sigma_j(\tau, p_j) = \sigma_j(\tau, p_{\text{ref}}) - N_j \ln \frac{p_j}{p_{\text{ref}}}$$

recovering Eq. (2.17a) of Turns (2000) [13].

## 17. ENTHALPY

17.1. **Stoichiometry.** cf. pp. 18, Chapter 2 Combustion and Thermochemistry, Section Reactant and Product Mixtures, Subsection Stoichiometry of Turns (2000) [13].

(53) 
$$C_x H_y + a(O_2 + 3.76 N_2) \longrightarrow x CO_2 + \left(\frac{y}{2}\right) H_2 O + 3.76 a N_2$$

and so  $a = x + \frac{y}{4}$ .

Therefore, the *stoichiometric air-fuel ratio* in this case is

$$\left(\frac{A}{F}\right)_{\text{stoic}} \equiv \left(\frac{m_{\text{air}}}{m_{\text{fuel}}}\right)_{\text{stoic}} = \frac{4.76 a \, MW_{\text{air}}}{1 \, MW_{\text{fuel}}}$$

where  $MW_{\text{air}}$ ;  $MW_{\text{fuel}}$  are molecular weights of air and fuel.

To reproduce Table 2.1 on pp. 19 of Turns (2000) [13], consider the following: reusing Eq. [53](#), then for methane,  $\text{CH}_4$ ,  $x = 1, y = 4$  in this case, and so  $a = 2$ . For  $\text{H}_2 + \text{O}_2$ , consider  $2 \text{H}_2 + \text{O}_2 \longrightarrow 2 \text{H}_2\text{O}$ . For  $\text{C(s)}^+$  air, consider a “1 to 1” reaction, of 1 C and 1 “air” reacting with each other.

Then consider how the  $O/F$  ratio is **defined**:

(54) 
$$\frac{O}{F} := \frac{m_{\text{Ox}}}{m_F} = \frac{\nu_{\text{Ox}} MW_{\text{Ox}}}{MW_F} \equiv \frac{\nu_{\text{Ox}} \widehat{m}_{\text{Ox}}}{\nu_F \widehat{m}_F}$$

Thus, for example, for  $\text{CH}_4$  and air,  $\nu_{\text{CH}_4} = 1, \nu_{\text{air}}$ , for  $\text{H}_2 + \text{O}_2$ ,  $\nu_{\text{H}_2} = 2$  and  $\nu_{\text{O}_2} = 1$ , and  $\nu_{\text{C(s)}} = 1, \nu_{\text{air}} = 1$ . Then look up the molecular masses or “molecular weights” for the constituents.

I suggest using the Python library **thermopy3** in this [github repository](#) in the following manner (and the following code snippet is in **combustion.py** of my github repository for **Propulsion**:

```
import thermopy3
from thermopy3 import nasa9polynomials as nasa9p
nasa9pDB = nasa9p.Database() # this initializes or ‘‘instantiates’’ the database that contains the
# coefficients and other data using NASA 9 polynomials
```

```
# find if a compound is in the database using the list_compound method
# then initialize or ‘‘create’’ the compound by the method set_compound
nasa9pDB.list_compound("methane")
CH4 = nasa9pDB.set_compound("methane")
```

```
nasa9pDB.list_compound("N2")
N2 = nasa9pDB.set_compound('N2')
```

```
nasa9pDB.list_compound("O2")
```

```
O2 = nasa9pDB.set_compound("O2")
```

```
nasa9pDB.list_compound("H2")
H2 = nasa9pDB.set_compound("H2")
```

```
nasa9pDB.list_compound("C")
C = nasa9pDB.set_compound("C")
```

```
CH4.molecular_weight # 0.016042459999999998 kg/mol
O2.molecular_weight # 0.0319988 kg/mol
N2.molecular_weight # 0.0280134 kg/mol
airMW = O2.molecular_weight + 3.76*N2.molecular_weight # 0.137329184 kg/mol
```

```
# O/F for oxidizer air and fuel methane
2*airMW / CH4.molecular_weight # 17.12071390547335
```

```
# O/F for oxidizer oxygen and fuel hydrogen
O2.molecular_weight / (2.*H2.molecular_weight) # 7.936682739051927
```

```
# O/F for oxidizer air and fuel Carbon
airMW / C.molecular_weight # 11.433903436102808
```

*equivalence ratio*  $\Phi$

$$\Phi := \frac{(A/F)_{\text{stoic}}}{(A/F)} = \frac{(F/A)}{(F/A)_{\text{stoic}}}$$

fuel rich  $\Phi > 1$

lean  $\Phi < 1$

So the immediate previous expressions are in Turns’ notation[13]. In the notation used above, the stoichiometric oxidizer-fuel ratio is

$$\left(\frac{m_{\text{Ox}}}{m_{\text{fuel}}}\right)_{\text{stoic}} = \frac{N_{\text{Ox}} \widehat{m}_{\text{Ox}}}{N_{\text{fuel}} \widehat{m}_{\text{fuel}}} = \frac{n_{\text{Ox}} \widehat{m}_{\text{Ox}}}{n_{\text{fuel}} \widehat{m}_{\text{fuel}}}$$

17.2. **Absolute (or standardized) Enthalpy and Enthalpy of Formation.** cf. Turns (2000) [13]

*enthalpy of formation*  $h_f$  - takes into account energy associated iwth chemical bonds

*sensible enthalpy change*  $\delta h_s$  - enthalpy associated only with temperature

absolute (or standard) enthalpy  $\bar{h}_i(\tau) := \bar{h}_{f,i}^0(\tau) + \Delta \bar{h}_{s,i}(\tau_{\text{ref}})$  where

$\bar{h}_i(\tau) \equiv$  absolute enthalpy at temperature  $\tau$

$\bar{h}_{f,i}^0(\tau_{\text{ref}}) \equiv$  enthalpy of formation at standard reference state  $(\tau_{\text{ref}}, p^0)$

$\Delta \bar{h}_{s,i}(\tau_{\text{ref}}) \equiv$  sensible enthalpy change in going from  $\tau_{\text{ref}}$  to  $\tau$

$p_{\text{ref}} = p^0 = 1 \text{ atm}$

Example 2.3.[13] gas stream at 1 atm. a mixture of  $\text{CO}$ ,  $\text{CO}_2$ ,  $\text{N}_2$ ; given  $X_{\text{CO}} = 0.1$ , gas stream temperature  $T = 1200 \text{ K}$

$$X_{\text{CO}_2}$$

$$X_{\text{N}_2} = 1 - X_{\text{CO}} - X_{\text{CO}_2} = 0.7$$

## 18. THERMOCHEMISTRY OF COMBUSTION

cf. 20151203 Dr. Polk AE121a Fall 2015 Lecture

The key is how we define enthalpy. We need to agree on

- (1) how to include enthalpy due to chemical changes in addition to the sensible enthalpy (due to temperature changes)
- (2) Reference state for sensible and chemical enthalpy

Taking into account the above, then

$$H_P(T_P) = \sum_i^{\text{products}} \mathcal{N}_i \left[ \underbrace{\Delta_f \bar{h}_i^0}_{\substack{\text{chemical} \\ \text{(enthalpy of formation)}}} + \underbrace{\bar{h}_i(T_P) - \bar{h}_i^0}_{\text{sensible}} \right]$$
$$H_R(T_R) = \sum_i^{\text{reactants}} \mathcal{N}_i \left[ \Delta_f \bar{h}_i^0 + \bar{h}_i(T_R) - \bar{h}_i^0 \right]$$

where  $P$  denotes products and  $R$  denotes reactants, and notation is Polk’s.

To reiterate, in my notation,

$$H_P(\tau_P) = \sum_{i \in \{\text{products}\}} n_i \left[ \Delta_f \bar{h}_i^0 + \bar{h}_i(\tau_P) - \bar{h}_i^0 \right]$$
$$H_R(\tau_R) = \sum_{i \in \{\text{reactants}\}} n_i \left[ \Delta_f \bar{h}_i^0 + \bar{h}_i(\tau_R) - \bar{h}_i^0 \right]$$

(55)

18.1. **Law of Mass Action.** Let’s recall the *law of mass action*, derived from Kittel and Kroemer [8] and mentioned for how chemical reactions go in Ch. 3 Chemical Rockets pp. 77 of Oates (1997) [1].

For an ideal gas,  $pV = N\tau$  or  $p = n\tau$ ,

$$\prod_j n_j^{\nu_j} = K(\tau) := \prod_j n_{Q_j}^{\nu_j} \exp(-\nu_j F_j(\text{int})/\tau) \xrightarrow{\tau^{\sum_j \nu_j}} \prod_j p_j^{\nu_j} \equiv K_p(\tau) \text{ equilibrium constant}$$
$$\prod_j X_j^{\nu_j} = K_p(\tau) p^{-(\nu_1 + \nu_2 + \cdots + \nu_n)}$$

(56)

where  $\frac{p_j}{p} = \frac{n_j}{\sum_i n_i} = X_j$  is the usual mole (molar; particle) fraction for species  $j$ .

18.1.1. *Example - H-O Reaction.* Consider



H :  $m + n = 1$   
O :  $m + 2q = 2l$   
cf. Chapter 7 Simplified Conservation Equations for Reacting Flows of Turns (2011) [13]  
Start from

$$m_i = \int_V \rho_i \text{vol}^n$$

(57)

and so

$$\frac{dm_i}{dt} \equiv \dot{m}_i = \int_V \frac{\partial \rho_i}{\partial t} \text{vol}^n + \int_{\partial V} \rho_i i_{u_i} \text{vol}^n = \int_V \frac{\partial \rho_i}{\partial t} \text{vol}^n + \int_{\partial V} \rho_i u_i^j dS_j = \dot{m}_i = \int_V \dot{m}_i''' \text{vol}^n$$

where we’ve defined

$$\dot{m}_i''' \equiv \text{rate of mass production per unit volume}$$

Thus

$$\implies \frac{\partial \rho_i}{\partial t} + \text{div}(\rho_i u_i) = \dot{m}_i'''$$

where

$$\text{div}(\rho_i u_i) := \frac{1}{\sqrt{g}} \frac{\partial(\rho_i u^j \sqrt{g})}{\partial x^j}$$

Now

$$\rho_i = \frac{M_i N_i}{V}$$
$$\implies \sum_i \rho_i = \rho \equiv \frac{M}{V} \text{ with } \sum_i M_i N_i \equiv M$$

where we assume negligible interaction between species.

Now

$$\sum_i \rho_i \mathbf{u}_i = \rho \mathbf{u} \text{ since}$$
$$\sum_i \rho_i \mathbf{u}_i = \sum_i \frac{M_i N_i \mathbf{u}_i}{V} = \frac{M \mathbf{u}}{V} = \rho \mathbf{u}$$

and so for

$$Y_i := \frac{m_i}{m}$$

Then

$$\mathbf{u} = \sum_i Y_i \mathbf{u}_i = \sum_i \frac{m_i}{m} \mathbf{u}_i = \sum_i \frac{M_i N_i \mathbf{u}_i}{M}$$

(58)

is the **bulk velocity** in usual fluid equations.

So taking the sum over species labeled by  $i$ ,  $\sum_i$ ,

$$\sum_i \frac{\partial \rho_i}{\partial t} + \text{div}(\sum_i \rho_i u_i) = \sum_i \dot{m}_i''' \implies \frac{\partial \rho}{\partial t} + \text{div}(\rho u) = \dot{m}''' = 0$$

(59)

Use Galilean invariant, or rather *Galilean transformation*, between “lab frame” (denoted with unprimed notation), and fluid-at-rest frame, or bulk velocity-is-zero frame (primed notation).

$$\mathbf{u}'_i = \mathbf{u}_i - \mathbf{u} \text{ so } \mathbf{u}_i = \mathbf{u}'_i + \mathbf{u}$$

(60)

Thus

$$\frac{\partial \rho_i}{\partial t} + \text{div}(\rho_i (u'_i + u)) = \frac{\partial \rho_i}{\partial t} + \text{div}(\rho_i u'_i) + \text{div}(\rho_i u) = \dot{m}_i'''$$

(61)

where

$$\begin{array}{ll} \text{div}(\rho_i u'_i) & \text{diffusion} \\ \text{div}(\rho_i u) & \text{convection} \end{array}$$

(62)

Define **diffusion flux**

$$\mathbf{j}_i = \rho_i \mathbf{u}'_i$$

(63)

Now, note that

$$Y_i = \frac{m_i}{m} = \frac{M_i N_i}{\sum_j M_j N_j} \text{ for}$$
$$Y_i m = m_i$$
$$Y_i \rho = \rho_i$$

and so

$$\frac{\partial(\rho Y_i)}{\partial t} + \text{div}(\rho Y_i u'_i) + \text{div}(\rho Y_i u) = \dot{m}_i'''$$
$$\implies \rho \frac{\partial Y_i}{\partial t} + Y_i \frac{\partial \rho}{\partial t} + Y_i \text{div}(\rho u) + \rho u^j \frac{\partial Y_i}{\partial x^j} = -\text{div} j_i + \dot{m}_i'''$$
$$\xrightarrow{\frac{\partial \rho}{\partial t} + \text{div}(\rho u) = 0} \rho \frac{\partial Y_i}{\partial t} + \rho u^j \frac{\partial Y_i}{\partial x^j} + \text{div} j_i = \dot{m}_i'''$$

To recap our results, armed with the definitions (and resulting identities)  $\rho_i := \frac{M_i N_i}{V} \equiv M_i n_i$ ,  $Y_i := \frac{m_i}{m} = \frac{\rho_i}{\rho}$  and so  $Y_i \rho = \rho_i$ , where the identity

and starting from  $m_i = \int_V \rho_i \text{vol}^n$ , and using Galilean transformation  $u_i = u'_i + u$

$$\begin{aligned} \dot{m}_i &= \int_V \left( \frac{\partial \rho_i}{\partial t} + \text{div}(\rho_i u_i) \right) \text{vol}^n = \int_V \text{vol}^n \left[ \frac{\partial \rho Y_i}{\partial t} + \text{div}(\rho_i (u'_i + u)) \right] = \int_V \text{vol}^n \left[ Y_i \frac{\partial \rho}{\partial t} + \rho \frac{\partial Y_i}{\partial t} + \text{div}(\rho_i u'_i) + \text{div}(\rho Y_i u) \right] = \\ &= \int_V \text{vol}^n \left[ Y_i \frac{\partial \rho}{\partial t} + \rho \frac{\partial Y_i}{\partial t} + \text{div}(\rho_i u'_i) + Y_i \text{div}(\rho u) + \rho u^j \frac{\partial Y_i}{\partial x^j} \right] = \int_V \text{vol}^n \left[ \rho \frac{\partial Y_i}{\partial t} + \text{div}(\rho_i u'_i) + \rho u^j \frac{\partial Y_i}{\partial x^j} \right] \end{aligned}$$

where total mass conservation  $\dot{m} = 0 \implies \frac{\partial \rho}{\partial t} + \text{div}(\rho u) = 0$  was used in the last equality, which must be true, even if chemical reactions occur.

Then use Fick's law for particle diffusion for each species:

$$\rho_i u'_i = j_{m_i} = -\rho D_{AB} \text{grad} Y_i$$

and so

$$(64) \quad \boxed{\dot{m}_i = \int_V \text{vol}^n \left[ \rho \frac{\partial Y_i}{\partial t} + -\text{div}(\rho D_{AB} \text{grad} Y_i) + \rho u^j \frac{\partial Y_i}{\partial x^j} \right]}$$

Compare this to Eq. 3.31, where steady flow is assumed ( $\partial(\rho Y_A)/\partial t = 0$ ) in the section on **Species Conservation** of Chapter 3 Introduction to Mass Transfer of Turns (2011) [13]. I don't agree with his usage of this steady flow assumption because one needs to use the *total* mass conservation first as shown above. Eq. 3.31 of Turns (2011) [13] is, for reference,

$$\dot{m}_A''' - \frac{d}{dx} \left[ Y_A \dot{m}'' - \rho D_{AB} \frac{dY_A}{dx} \right] = 0$$

What is useful from Turns (2011) [13] is the physical interpretation he provides: from Eq. 64,  $\dot{m}_i$  is the net rate of production of species  $i$  by chemical reaction by the entire system.  $\rho \frac{\partial Y_i}{\partial t} + -\text{div}(\rho D_{AB} \text{grad} Y_i) + \rho u^j \frac{\partial Y_i}{\partial x^j}$  is the net flow of species  $i$  out of the (control) volume  $V$  per unit volume.

Then take a look at the section on **Species Mass Conservation (Species Continuity)** on pp. 218 of Chapter 7 “Simplified Conservation Equations” of Turns (2011) [13], namely Eq. (7.8) and (7.10). Further physical interpretation is provided there, and so we can interpret terms in Eq. 64 as follows:

$-\text{div}(\rho D_{AB} \text{grad} Y_i)$	mass flow of species $i$ due to molecular diffusion per unit volume ( $\text{kg}/s \cdot m^3$ )
$\rho u^j \frac{\partial Y_i}{\partial x^j}$	mass flow of species $i$ due to convection (advection by bulk flow) per unit volume ( $\text{kg}/s \cdot m^3$ )
$\dot{m}_i$	net mass production rate of species $i$ by chemical reaction (for the entire system) $\text{kg}/s$

**18.2. Energy Conservation applied to Combustion.** Following the Section **Energy Conservation** on pp. 334, Chapter 7 Simplified Conservation Equations, of Turns (2011) [13],

$$\dot{E} = - \int_V \text{div}((k + h')\mathbf{u}) \text{vol}^n$$

where we keep in mind the difference in notation and meaning of  $h$  and  $h'$ :  $h' := \frac{H}{V}$  where  $H = \int_V h' \text{vol}^n$ , is the enthalpy density, whereas  $h := \frac{H}{m}$  where  $H = \int_V h m = \int_V h \rho \text{vol}^n$  is the so-called *specific enthalpy*, enthalpy per unit mass. So  $\rho h = h'$

Noting that

$$\text{div}\left(\left(\frac{\rho u^2}{2} + h'\right)u\right) = \text{div}\left(\left(\frac{\rho u^2}{2} + \rho h\right)u\right) = \text{div}(\rho u)\left(h + \frac{u^2}{2}\right) + \rho u^j \frac{\partial}{\partial x^j} \left(h + \frac{u^2}{2}\right)$$

and that

$$\begin{aligned} u^2 = g_{ij} u^i u^j \implies u^k \frac{\partial}{\partial x^k} u^2 &= 2u^k g_{ij} \frac{\partial u^i}{\partial x^k} u^j + \frac{\partial g_{ij}}{\partial x^k} u^k u^i u^j = 2u^k \langle u, \frac{\partial u}{\partial x^k} \rangle + 2g_{ks} \Gamma_{ij}^k u^i u^j u^s = \\ &= 2g_{ks} u^s \left[ u^j \frac{\partial u^k}{\partial x^j} + \Gamma_{ij}^k u^i u^j \right] = 2g_{ks} u^s (\nabla_u u) = 2\langle u, \nabla_u u \rangle \end{aligned}$$

$$\frac{\partial g_{ij}}{\partial x^k} \dot{\phi}^k \dot{\phi}^i \dot{\phi}^j = 2g_{ks} \Gamma_{ij}^k \dot{\phi}^i \dot{\phi}^j \dot{\phi}^s \text{ with } u^i = \dot{\phi}^i$$

was used, which is in the proof for Theorem 3.25 on pp. 49 of section 3.8 The natural Lagrangian on manifolds of Calin and Chang (2005) [17]

Then

$$\dot{E} = - \int_V \text{vol}^n \left[ \text{div}(\rho u) \left(h + \frac{u^2}{2}\right) + \rho u^j \frac{\partial}{\partial x^j} \left(h + \frac{u^2}{2}\right) \right] = \int_V \frac{\partial \rho}{\partial t} \left(h + \frac{u^2}{2}\right) \text{vol}^n - \int_V \text{vol}^n \left[ \rho u^j \frac{\partial h}{\partial x^j} + \rho \langle u, \nabla_u u \rangle \right]$$

Now consider the heat  $Q \in \Omega^1(\Sigma)$ . If we define the heat density  $q' \in C^\infty(\mathbb{R} \times N) \times \Omega^1(\Sigma)$ , then

$$\dot{Q} = \int_V \left( \frac{\partial q'}{\partial t} + \text{div}(q' u) \right) \text{vol}^n \xrightarrow{\frac{\partial q'}{\partial t} = 0} \int_V \text{vol}^n \left[ u^j \frac{\partial q'}{\partial x^j} + q' \text{div} u \right]$$

where we'd have to deal with whether the fluid is compressible  $\text{div} u \neq 0$  or incompressible  $\text{div} u = 0$ .

If we define the “specific heat,” the heat per unit mass,  $q$ , defined via  $Q = \int \rho q \text{vol}^n$ ,  $q \in C^\infty(\mathbb{R} \times N) \times \Omega^1(\Sigma)$ ,

$$\begin{aligned} \dot{Q} &= \int_V \left( \rho \frac{\partial q}{\partial t} + q \frac{\partial \rho}{\partial t} \right) \text{vol}^n + \text{div}(\rho q u) \text{vol}^n = \int_V \left( \rho \frac{\partial q}{\partial t} + q \frac{\partial \rho}{\partial t} \right) \text{vol}^n + q \text{div}(\rho u) \text{vol}^n + \rho u^j \frac{\partial q}{\partial x^j} \text{vol}^n = \\ &= \int_V \rho \text{vol}^n \left( \frac{\partial q}{\partial t} + u^j \frac{\partial q}{\partial x^j} \right) \end{aligned}$$

where total mass conservation was used (which should still be valid even if chemical reactions occur; the masses for individual species can change, but the total mass doesn't change for chemical reactions), so that

$$u^j \frac{\partial q}{\partial x^j} = -u^j \frac{\partial h}{\partial x^j} - \langle u, \nabla_u u \rangle \text{ or } -\frac{\partial q}{\partial x^j} = \frac{\partial h}{\partial x^j} + (\nabla_u u)_j$$

Compare this expression we had just derived

$$(65) \quad -\frac{\partial q}{\partial x^j} = \frac{\partial h}{\partial x^j} + (\nabla_u u)_j$$

with Eq. 7.51 on pp. 235 of Chapter 7 of Turns (2011) [13]

$$-\frac{d\dot{Q}_x''}{dx} = \dot{m}'' \left( \frac{dh}{dx} + v_x \frac{dv_x}{dx} \right)$$

## 19. DROPLET EVAPORATION

**19.1. Evaporation Rate.** cf. Chapter 3 Introduction to Mass Transfer, Droplet Evaporation section, Evaporation Rate Sub-Section of Turns (2011) [13].

Starting from  $m = \int_V \rho \text{vol}^n$ ,

$$(66) \quad \begin{aligned} \dot{m} &= \int_V \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial V} d(\rho i_u \text{vol}^n) = \int_V \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial V} \rho u^i dS_i \\ &= 0 + \rho u^r(r_s) 4\pi r_s^2 \equiv 4\pi r^2 \dot{m}'' \end{aligned}$$

where we had assumed, first, steady state for the density  $\rho$ , and second, spherical symmetry.

Consider  $m = \sum_i m_i$ . For binary collisions condition,  $m = m_A + m_B$ . Then

$\dot{m} = \sum_i \dot{m}_i$ . Assuming the above just like in Eq. 66, steady-state density and spherical symmetry, then  $\dot{m}'' = \dot{m}''_A + \dot{m}''_B = \dot{m}''_A$ , since  $\dot{m}''_B = 0$  (assumed *no flow into droplet*).

19.1.1. *Separating diffusion and bulk flow.* Beginning with  $m_i = \int_V \rho_i \text{vol}^n$  (and  $m = \sum_i m_i = \int_V \rho \text{vol}^n$ , where  $\rho = \sum_i \rho_i$ , which is allowed if we assume *negligible interaction* between species, until chemical reactions),

$$\dot{m}_i = \int_V \frac{\partial \rho_i}{\partial t} \text{vol}^n + \int_{\partial V} \rho_i u_i^j dS_j = \int_V \frac{\partial \rho_i}{\partial t} \text{vol}^n + \int_{\partial V} \rho_i (u_i')^j dS_j + \int_{\partial V} \rho_i u^j dS_j = \int_V \frac{\partial \rho_i}{\partial t} \text{vol}^n + \int_{\partial V} j_{i,\text{diff}} \cdot dS + \int_{\partial V} \rho_i u \cdot dS \text{ using}$$

$$\int_{\partial V} \rho_i u_i^j dS_j = \int_{\partial V} \rho_i ((u_i')^j + u^j) dS_j$$

Assuming steady state density  $\rho_i$ ,  
constant physical parameters on a particular choice of the surface boundary  $\partial V$ , with surface area  $A$ ,

$$\frac{\dot{m}_i}{A} \equiv \dot{m}_i'' = 0 + (j_{i,\text{diff}})^r + \rho_i u^r$$

Now for the total mass of all (i.e. all the species),  $m = \sum_i m_i$ ,

$$\dot{m} = \int \frac{\partial \rho}{\partial t} \text{vol}^n + \int_{\partial V} \rho u^j dS_j$$

Assuming steady state density  $\rho$ ,  
constant physical parameters on a particular choice of the surface boundary  $\partial V$ , with surface area  $A$ ,

$$\dot{m}/A \equiv \dot{m}'' = 0 + \rho u^r$$

So one *must define carefully* the choice of volume  $V$  (and thus surface  $\partial V$ ). For we want to consider a choice of volume  $V$  s.t. the choice of  $\partial V$  results in a *net flux* of the bulk flow (we want a measure of that across the surface  $\partial V$ . In this case,

$$\dot{m}'' = \rho u^r \equiv \rho u^n$$

From the definition of  $Y_i := \frac{m_i}{m}$ , then

$$\rho_i u = Y_i \rho u \equiv Y_i (\dot{m}'')$$

Thus

$$(67) \quad \dot{m}_i'' = Y_i (\dot{m}'') + (j_{i,\text{diff}})$$

which is Eqns. (3.1), (3.5) in Turns (2011) [13].

From Le Bellac, Mortessagne, Batrouni (2004) [15], namely Chapter 6 on Irreversible Processes, for the section on “Particle Diffusion”, for particle diffusion of species  $A$ , with the number of particles of species  $A$ ,

$$j_\alpha^{NA} = \sum_{j,\beta} L_{NAj}^{\alpha\beta} \partial_\beta \gamma_j = \sum_{j,\beta} L_{NAj}^{\alpha\beta} \partial_\beta \left( \frac{-\mu_j}{\tau} \right) = \sum_{j,\beta} \frac{-L_{NAj}^{\alpha\beta}}{\tau} \sum_i \frac{\partial \mu_j}{\partial n_i} \frac{\partial n_i}{\partial x^\beta}$$

If  $\mu_i = \mu_i(n_i)$ , so that  $\mu_i$  doesn’t depend on  $\mu_j$ ,  $j \neq i$ , i.e. the chemical potential of species  $i$  does not depend on the concentrations of the other species,

$$\sum_{j,\beta} \frac{-L_{NAj}^{\alpha\beta}}{\tau} \frac{\partial \mu_j}{\partial n_j} \frac{\partial n_j}{\partial x^\beta} = \sum_{j,\beta} \frac{-L_{NAj}^{\alpha\beta}}{\tau} \frac{1}{\kappa_T n_j^2} \frac{\partial n_j}{\partial x^\beta}$$

and so

$$(68) \quad \mathbf{j}_{NA} = \sum_j \frac{-L_{NAj}}{\tau \kappa_T n_j^2} \nabla n_j$$

Now

$$\frac{N_i}{V} \equiv n_i = \frac{Y_i m / M_i}{V} = \frac{Y_i m}{M_i V} = \frac{Y_i \rho}{M_i} \text{ so } \boxed{\nabla n_i = \nabla \frac{Y_i \rho}{M_i} = \frac{\rho}{M_i} \nabla Y_i} \text{ for}$$

$$\sum_i \rho_i = \rho := \frac{M}{V} \implies \sum_i M_i n_i = \frac{M}{V}$$

So then we can reexpress Eqn. 68 in terms of *mass fractions*  $Y_i$  as

$$\mathbf{j}_{NA} = -\rho \sum_j \frac{L_{NAj}}{\tau \kappa_T n_j^2 M_j} \nabla Y_j$$

Note that  $\kappa_T$  is the **coefficient of isothermal compressibility** which was defined in Le Bellac, Mortessagne, Batrouni (2004) [15],

$$\kappa_T = \frac{-1}{V} \left. \frac{\partial V}{\partial p} \right|_\tau$$

Now I guess that  $L_{NAj} \propto \delta_{NAj}$  meaning  $L_{NAj} \neq 0$  only if  $j = N_A$ , and so  $\mathbf{j}_{NA}$  only depends on the gradient  $\nabla Y_{N_A}$ , the gradient for what is the concentration of species  $A$ , and *not* on the gradient of the concentrations of the other species. I *think* (please contact me if this is wrong) this is because the species are nonreacting, and even when including chemical reactions, the chemical reactions does not affect this diffusion process; the diffusion process is independent of what goes on with the chemical reactions after the chemical reactions occur.

I conjecture that  $L_{NAj}$  does not depend on direction (i.e. *isotropic*), which is reasonable, by spatial symmetry.

If so, then

$$(69) \quad \mathbf{j}_{NA} = -\rho \left( \frac{L_{N_A N_A}}{\tau \kappa_T n_A^2 M_A} \right) \nabla Y_A$$

19.1.2. *Particle Diffusion.* cf. pp. 399 of Kittel and Kroemer [8]

Consider a system.

One end in diffusive contact with reservoir at chemical potential  $\mu_1$

Other end in diffusive contact with reservoir at chemical potential  $\mu_2$

Constant temperature  $\tau$ .

If  $\mu_1 > \mu_2$ , particle flow through system from reservoir 1 to reservoir 2;  $1 \rightarrow 2$ .

$n_i \equiv$  particle concentration in  $i$

Take

$$(70) \quad \mathbf{j}_n = -D \text{grad} n$$

which is **Fick’s law**, and where  $D \equiv$  particle diffusion constant or **diffusivity**.

Mean free path  $l$ . Particles freely travel over  $l$ .

Assume in a collision at  $z$ , particles come into local equilibrium at local chemical potential  $\mu(z)$ , local concentration  $n(z)$ .

At  $z$ , particle flux density in positive  $z$  direction  $\frac{1}{2} n(z - l_z) \bar{c}_z$

particle flux density in negative  $z$  direction  $-\frac{1}{2} n(z + l_z) \bar{c}_z$

Note  $n(z - l_z)$  is particle concentration at  $z - l_z$

$$J_n^z = \frac{1}{2} [n(z - l_z) - n(z + l_z)] \bar{c}_z = -\frac{dn}{dz} \bar{c}_z l_z$$

where  $\bar{c}_z = \bar{c} \cos \theta$

$$\bar{l}_z = \bar{l} \cos \theta$$

$$\langle \bar{c}_z l_z \rangle = \bar{c} \bar{l} \frac{\int_{\text{hemisphere}} \cos^2 \theta dS}{\int_{\text{hemisphere}} dS} = \bar{c} \bar{l} \frac{\int_0^{\pi/2} d\theta \int_0^{2\pi} d\varphi \cos^2 \theta \sin \theta d\theta}{\int_0^{\pi/2} d\theta \sin \theta \int_0^{2\pi} d\varphi} = \bar{c} \bar{l} \frac{1}{3}$$

Comparing with Fick’s law,

$$J_n^z = \frac{-1}{3} \bar{c} \bar{l} \frac{dn}{dz} \text{ or } \mathbf{J}_N = -\frac{1}{3} \bar{c} \bar{l} \nabla n$$

For diffusivity  $D$  is then  $D = \frac{1}{3} \bar{c} \bar{l}$ .



Now recall that  $\bar{l}$ , the *mean free path*, was derived from kinetic theory:

$$l = \frac{1}{n\pi d^2}$$

where  $d$  is the diameter of the particle.

The mean thermal velocity was derived from the Maxwell distribution:

$$\bar{c} = \left( \frac{8\tau}{M\pi} \right)^{1/2}$$

**19.2. Binary case for Fick’s law for mass transfer; A and B species only.** It seems that for combustion, we only worry about 2 species,  $A$  and  $B$ , coming together, most of the time.

I am following pp. 84 Chapter 3 Introduction to Mass Transfer, Section **Mass Transfer Rate Laws** of Turns (2011) [13] here, applying Fick’s law of Diffusion from my own development and comparing results from Turns.

Recall the basic quantities for summing up masses and particles (in this case, only two kinds), and mass fraction vs. mole fraction:

$$\begin{aligned} N_A + N_B &= N \\ m_A + m_B &= m \\ Y_A &:= \frac{m_A}{m} = \frac{M_A N_A}{M_A N_A + M_B N_B} \\ m_A + m_B &= Y_A m + Y_B m = m \\ X_A &:= \frac{N_A}{N} \\ M_A X_A + M_B X_B &= \frac{m}{N} \\ \rho_A &= \frac{M_A N_A}{V} \text{ so } Y_A(\rho_A + \rho_B) = \rho_A \end{aligned}$$

Now

$$\dot{m}_A = \frac{d}{dt} \int_V \rho_A \text{vol}^n = \int_V \left[ \frac{\partial \rho_A}{\partial t} + \text{div}(\rho_A u_a) \right] \text{vol}^n$$

Using Galilean transformation  $\mathbf{u}_A - \mathbf{u} = \mathbf{u}'_A$  or  $\mathbf{u}_A = \mathbf{u} + \mathbf{u}'_A$ ,

$$\text{div}(\rho_A u) = \text{div}(\rho_A (u + u'_A)) = \text{div}(\rho_A u) + \text{div}(\rho_A u'_A)$$

Considering the term  $\text{div}(\rho_A u)$ ,

$$\text{div}(\rho_A u) = \text{div}(Y_A \rho u) \text{ and also note that}$$

$$\int_V \text{div}(Y_A \rho u) = \int_{\partial V} dS \cdot Y_A \rho \mathbf{u}$$

Now for term  $\text{div}(\rho_A u'_A)$ , using Eq. 70, Fick’s law,

$$\text{div}(\rho_A u'_A) = \text{div}(j_{\rho_A}) = \text{div}(-M_A \mathcal{D}_{AB} \text{grad} n_A) = -\text{div}(\mathcal{D}_{AB} \text{grad} \rho_A) = -\text{div}(\rho \mathcal{D}_{AB} \text{grad} Y_A)$$

Putting everything in terms of a surface area integral,

$$\dot{m}_A = \int_{\partial V} dS \cdot \dot{m}_A'' = \int_{\partial V} dS \cdot \rho u Y_A - \int_{\partial V} dS \cdot \rho \mathcal{D}_{AB} \text{grad} Y_A$$

Then

$$(71) \quad \boxed{\dot{m}_A'' = Y_A \rho u - \rho \mathcal{D}_{AB} \text{grad} Y_A}$$

Note that  $\mathcal{D}_{AB}$  is in units of  $[m^2/2]$ .

The physical interpretation for the terms is as follows:

$\dot{m}_A''$  is the mass flow of species  $A$  per unit area

$Y_A \rho u$  is the mass flow of species  $A$  associated with bulk flow per unit area

$\rho \mathcal{D}_{AB} \text{grad} Y_A$  is the mass flow of species  $A$  associated with molecular diffusion per unit area.

Compare this to Eq. 3.1, 3.5 of Turns (2011) [13]:

$$(72) \quad \dot{m}_A'' = Y_A (\dot{m}_A'' + \dot{m}_B'') - \rho \mathcal{D}_{AB} \text{grad} Y_A$$

## 20. DROPLET EVAPORATION AND BURNING

I am following Chapter 10 Droplet Evaporation and Burning of Turns (2011) [13].

“There are two types of liquid rockets: **pressure-fed**, in which the fuel and oxidizer are pushed into the combustion chamber by a high-pressure gas; and **pump-fed**, where turbopumps deliver the propellants.” [13]

## 21. SHVAB-ZELDOVICH FORMS

Lewis number ( $\text{Le} := \frac{K}{\rho c_p D}$ )  $\text{Le} = 1$  by assumption.

Species flux and Fick’s law.

$$\mathbf{j}_q = -K \text{grad} \tau + \sum_i \dot{m}_i''_{\text{diff}} h_i = -K \text{grad} \tau - \sum_i \rho D (\text{grad} Y_i) h_i = -K \text{grad} \tau - \rho D \text{grad} \sum_i Y_i h_i + \rho D \sum_i Y_i \text{grad} h_i$$

Noting that

$$\sum_i Y_i \text{grad} h_i = \sum_i Y_i c_{pi} \text{grad} \tau = c_p \text{grad} \tau$$

then

$$j_q = -K \text{grad} \tau - \rho D \text{grad} h + \rho D c_p \text{grad} \tau$$

thermal diffusivity  $\alpha := \frac{K}{\rho c_p}$ .

If  $\text{Le} = 1$ ,  $K = \rho D c_p$  and so  $j_q = -\rho D \text{grad} h$  or  $\dot{Q}'' = -\rho D \text{grad} h$ .

**21.1. Heat  $Q$ .** Again  $Q \in \Omega^1(\Sigma)$  represents heat on the system.

The conundrum is whether  $Q$  is best served with the heat density  $q' := Q/V$  or “specific heat” per mass,  $q := Q/m$ , because it’ll affect whether we have to deal with the compressibility or incompressibility of the fluid or mass conservation, and the form of the “heat flux.”

So

$$Q = \int_V q \rho \text{vol}^n = \int_V q m = \int_V q' \text{vol}^n$$

Then

$$\begin{aligned} \dot{Q} &= \int_V \left( \rho \frac{\partial q}{\partial t} + q \frac{\partial \rho}{\partial t} \right) \text{vol}^n + \int_V \text{div}(\rho q u) \text{vol}^n = \int_V \left[ \rho \left( \frac{\partial q}{\partial t} + u^j \frac{\partial q}{\partial x^j} \right) + q \left( \frac{\partial \rho}{\partial t} + \text{div}(\rho u) \right) \right] \text{vol}^n = \\ &= \int_V \left( \frac{\partial q'}{\partial t} + \text{div}(q' u) \right) \text{vol}^n = \int_V \left( \frac{\partial q'}{\partial t} + \text{div}(\mathbf{j}_{q'}) \right) \text{vol}^n \end{aligned}$$

where

$$\text{div}(\mathbf{j}_{q'}) = \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^j} (-\rho D \sqrt{g} \text{grad} h) = \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^j} (-\rho D \frac{\partial h}{\partial x^k} g^{jk} \sqrt{g})$$

Employ the definition of absolute (or standardized) enthalpy in Turns (2011) [13]:

$$h = \sum_i Y_i \Delta_f \bar{h}_i^\circ + \int_{T_{\text{ref}}}^T c_p d\tau$$

The proper way to treat the usage of energy conservation with the Shvab-Zeldovich form is this: starting from  $Q = dE - W$ , and having no external work done on the system  $\dot{W} = 0$

$$(73) \quad Q \left( \frac{d}{dt} \right) = \dot{Q} = \int_V \text{vol}^n \left( \frac{\partial q'}{\partial t} + \text{div}(j_{q'}) \right) = dE \left( \frac{d}{dt} \right) = \dot{E} = - \int_V \text{vol}^n \left[ (\rho u^j) \left( \frac{\partial h}{\partial x^j} + (\nabla_u u)_j \right) \right]$$

where on the right hand side (RHS) of Eq. 73 is from the energy conservation of the total energy density of the fluid system,

$$\frac{\partial \epsilon}{\partial t} + \text{div}((h' + k)u) = 0$$



which is from 6.5.2 of Le Bellac, Mortessagne, Batrouni (2004) [15]). Then for steady-state assumption,  $\frac{\partial q'}{\partial t} = 0$ ,

$$\text{div}(j_{q'}) = -(\rho u^j) \left( \frac{\partial h}{\partial x^j} + (\nabla_u u)_j \right)$$

The Shvab-Zeldovich form for  $j_{q'}$  is

$$(74) \quad j_q = -\rho D \text{grad} h$$

and so

$$(75) \quad \text{div}(\mathbf{j}_{q'}) = \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^j} (-\rho D \sqrt{g} \text{grad} h) = \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^j} (-\rho D \frac{\partial h}{\partial x^k} g^{jk} \sqrt{g}) = -(\rho u^j) \left( \frac{\partial h}{\partial x^j} + (\nabla_u u)_j \right)$$

Compare Eq. 75 to Eq. (7.61) of Turns (2011) [13].

Employ now Turns' definition of absolute (or standardized) enthalpy [13]:

$$h = \sum_i Y_i \Delta_f \bar{h}_i^0 + \int_{T_{\text{ref}}}^T c_p dT$$

Then applying this definition to Eq. 75,

$$\begin{aligned} \text{div}(\rho D \text{grad} h) &= \text{div} \left( \rho D \sum_i \Delta_f \bar{h}_i^0 \text{grad} Y_i + \rho D \text{grad} \int_{T_{\text{ref}}}^T c_p dT \right) = \\ &= (\rho u^j) \left( \frac{\partial h}{\partial x^j} + (\nabla_u u)_j \right) = (\rho u^j) \left( \sum_i \Delta_f \bar{h}_i^0 \frac{\partial Y_i}{\partial x^j} + \frac{\partial}{\partial x^j} \int_{T_{\text{ref}}}^T c_p dT + (\nabla_u u)_j \right) \text{ or (moving terms from left to right and vice versa)} \\ &(\rho u^j) \frac{\partial}{\partial x^j} \int c_p dT - \text{div}(\rho D \text{grad} \int c_p dT) + \rho u^j (\nabla_u u)_j = \sum_i \text{div}(\rho D \Delta_f \bar{h}_i^0 \text{grad} Y_i) - (\rho u^j) \Delta_f \bar{h}_i^0 \frac{\partial Y_i}{\partial x^j} \end{aligned}$$

Use Eq. 64, the generalized equation for species conservation,

$$\sum_i \Delta_f \bar{h}_i^0 (\text{div}(\rho D \text{grad} Y_i) - (\rho u^j) \frac{\partial Y_i}{\partial x^j}) = \sum_i \Delta_f \bar{h}_i^0 (-\dot{m}_i''')$$

where  $\dot{m}_i = \int_V \text{vol}^n \dot{m}_i'''$ .

Thus, in general,

$$(76) \quad \boxed{(\rho u^j) \frac{\partial}{\partial x^j} \int c_p dT - \text{div}(\rho D \text{grad} \int c_p dT) + \rho u^j (\nabla_u u)_j = - \sum_i \Delta_f \bar{h}_i^0 \dot{m}_i'''}$$

where Turns (2011) [13] gives the (important) physical interpretations:

$(\rho u^j) \frac{\partial}{\partial x^j} \int c_p dT$	rate of sensible enthalpy transport by convection (advection) per unit volume ( $W/m^3$ )
$-\text{div}(\rho D \text{grad} \int c_p dT)$	rate of sensible enthalpy transport by diffusion per unit volume ( $W/m^3$ )
$-\sum_i \Delta_f \bar{h}_i^0 \dot{m}_i'''$	rate of formation enthalpy production by chemical reaction per unit volume ( $W/m^3$ )

Turns (2011) [13] says that  $-\sum_i \Delta_f \bar{h}_i^0 \dot{m}_i'''$  is the “rate of sensible enthalpy production by chemical reaction per unit volume ( $W/m^3$ )”, but I think this is wrong because the terms are dealing with the formation enthalpy, the enthalpy locked in the chemical bonds of the constituents.

It will pay to consider the case of spherical symmetry for droplets. Assuming only dependence upon  $r$  for all quantities by spherical symmetry, and keeping in mind that total mass conservation holds, i.e.

$$\begin{aligned} \dot{m} = 0 &\implies \frac{\partial \rho}{\partial t} + \text{div}(\rho u) = 0 \xrightarrow{\text{spherical symmetry}} \frac{1}{r^2} \frac{d}{dr} (\rho u^r r^2) = -\frac{\partial \rho}{\partial t} \\ &\xrightarrow{\frac{\partial \rho}{\partial t}=0 \text{ steady state}} \frac{1}{r^2} \frac{d}{dr} (\rho u^r r^2) = 0 \end{aligned}$$

then

$$\begin{aligned} -\text{div}(\rho D \text{grad} \int c_p dT) &= -\frac{1}{r^2} \frac{d(\rho D r^2 \frac{d \int c_p dT}{dr})}{dr} \\ \frac{d}{dr} (r^2 \rho u^r \int c_p dT) &= \left( \frac{d}{dr} (r^2 \rho u^r) \right) \int c_p dT + r^2 \rho u^r \frac{d}{dr} \int c_p dT \xrightarrow{\dot{m}=0} 0 + r^2 \rho u^r \frac{d}{dr} \int c_p dT \end{aligned}$$

so

$$(77) \quad \frac{1}{r^2} \frac{d}{dr} \left[ r^2 \rho u^r \int c_p dT - \rho D r^2 \frac{d \int c_p dT}{dr} \right] = - \sum_i \Delta_f h_i^0 \dot{m}_i'''$$

Compare this to Eq. (7.65) of Turns (2011) [13].

21.1.1. *Temperature distribution in the gas phase for Droplet Evaporation.* From Eq. 77, supposing steady state mass conservation with spherical symmetry:

$$\dot{m} = 4\pi r^2 \rho u^r$$

and supposing that

$$\sum_i \Delta_f \bar{h}_i^0 \dot{m}_i''' = 0$$

which means that “reaction rate term is zero, since no reactions occur for pure evaporation” (cf. pp. 372 Ch. 10 Droplet Evaporation and Burning, Eq. (10.3), of Turns (2011) [13]), and that the Lewis number  $Le := \frac{K}{\rho c_p D}$  is 1, then for the remaining terms in Eq. 77

$$\begin{aligned} r^2 \rho u^r \frac{d}{dr} \int c_p dT &= r^2 \rho u^r c_p \frac{dT}{dr} = \frac{\dot{m}}{4\pi} c_p \frac{dT}{dr} \\ \frac{d}{dr} \left[ \rho D r^2 \frac{d \int c_p dT}{dr} \right] &= \frac{d}{dr} \left[ \rho D r^2 c_p \frac{dT}{dr} \right] = \frac{d}{dr} \left[ K r^2 \frac{dT}{dr} \right] \\ &\implies \frac{d}{dr} \left[ r^2 \frac{dT}{dr} \right] = \frac{\dot{m} c_{pg}}{4\pi K} \frac{dT}{dr} \equiv Z \dot{m} \frac{dT}{dr} \end{aligned}$$

Let's solve this (easy) ODE (ordinary differential equation):

$$\frac{d(r^2 \frac{dT}{dr})}{dr} = Z \dot{m} \frac{dT}{dr}$$

Then

$$\begin{aligned} \implies 2r\dot{T} + r^2\ddot{T} &= Z\dot{m}\dot{T} \text{ or } \ddot{T} = \frac{(Z\dot{m} - 2r)}{r^2} \dot{T} \text{ so } \ln \dot{T} = \frac{Z\dot{m}}{-r} - 2 \ln r + C \text{ or } \dot{T} = \frac{C_1}{r^2} \exp \left( \frac{-Z\dot{m}}{r} \right) \\ &\xrightarrow{\int dr} T(r) = \frac{C_1 \exp \left( \frac{-Z\dot{m}}{r} \right)}{Z\dot{m}} + C_2 \end{aligned}$$

Let's consider the (important physically) boundary conditions:

$$(78) \quad \begin{aligned} T_b &\equiv T_{\text{boil}} = T(r = r_s) \\ T(r = \infty) &= T_\infty = \frac{C_1}{Z\dot{m}} + C_2 \end{aligned}$$

which is physically interpreted (importantly) that out in  $r = \infty$ , the system is in thermal equilibrium to some temperature  $T_\infty$ , and at the surface of the droplet,  $r = r_s$ , the liquid making up the droplet is evaporating away, and this occurs at the boiling temperature  $T_b$ . Thus, the temperature at  $r = r_s$ , the surface of the droplet, has to be  $T_b$ .

Then

$$\begin{aligned} T_\infty &= \frac{C_1}{Z\dot{m}} + C_2 & T_b &= \frac{C_1 \exp\left(-\frac{Z\dot{m}}{r_s}\right)}{Z\dot{m}} + C_2 \\ C_1 \left( \frac{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)}{Z\dot{m}} \right) &= T_\infty - T_b \\ \implies C_2 &= T_\infty - \frac{T_\infty - T_b}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)} = \frac{T_b - T_\infty \exp\left(-\frac{Z\dot{m}}{r_s}\right)}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)} \end{aligned}$$

Therefore

$$(79) \quad \boxed{T(r) = \frac{(T_\infty - T_b) \exp\left(-\frac{Z\dot{m}}{r}\right) + T_b - T_\infty \exp\left(-\frac{Z\dot{m}}{r_s}\right)}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)}}$$

Note that the dependence on  $r$  of  $T(r)$ ,  $\frac{dT(r)}{dr}$  is easily given by

$$(80) \quad \frac{dT(r)}{dr} = \frac{(T_\infty - T_b) \exp\left(-\frac{Z\dot{m}}{r}\right) (Z\dot{m} \frac{1}{r^2})}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)}$$

Let’s consider the surface of the droplet. Obviously, if it’s evaporating off the surface, then heat  $Q$  is supplied at the surface, raising the temperature of the droplet molecules at the surface to boiling temperature, enough to change (its state) into a vapor. This is the so-called latent heat of vaporization. Denote this heat by  $Q_{\text{conduction}}$ .

Starting from the general statement that  $Q = dH - Vdp$ , then the thermodynamic process  $\gamma : \mathbb{R} \rightarrow \Sigma$  we want to consider here results in the following:

$$\int Q(\dot{\gamma}) = H_{\text{vap}} - H_{\text{liq}}$$

Then

$$(81) \quad \dot{Q}_{\text{conduction}} = \frac{d}{dt} \int_\gamma Q(\dot{\gamma}) = \dot{m}(h_{\text{vap}} - h_{\text{liq}}) \equiv \dot{m}h_{fg}$$

Now consider Fourier’s law for heating (cf. pp. 401 Eq. (26) of Kittel and Kroemer (1980) [8]). It says

$$j_u = -K \text{grad} \tau$$

This implies, for the steady-state case, that

$$\dot{Q}_{\text{conduction}} = \int_{\partial V} dS_j j_u^j = K4\pi r_s^2 \left. \frac{dT}{dr} \right|_{r_s}$$

Note that the direction is such that positive heat flows *into* the droplet system, whereas, geometrically, the normal to the surface is outward and *out* of the droplet. Plugging this into Eq. 81,

$$\dot{m}h_{fg} = K4\pi r_s^2 \left. \frac{dT}{dr} \right|_{r_s}$$

Plugging in Eq. 80,

$$\begin{aligned} \frac{(T_\infty - T_b) \exp\left(-\frac{Z\dot{m}}{r_s}\right) (Z\dot{m} \frac{1}{r_s^2})}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)} \left( \frac{K4\pi r_s^2}{h_{fg}} \right) &= \dot{m} \implies \frac{(T_\infty - T_b) \exp\left(-\frac{Z\dot{m}}{r_s}\right)}{1 - \exp\left(-\frac{Z\dot{m}}{r_s}\right)} = \frac{h_{fg}}{c_{pg}} \\ &\implies \ln \left( \frac{c_{pg}}{h_{fg}} (T_\infty - T_b) + 1 \right) = \frac{Z\dot{m}}{r_s} \\ &\implies \dot{m} = \frac{4\pi K_g r_s}{c_{pg}} \ln \left( \frac{c_{pg}}{h_{fg}} (T_\infty - T_b) + 1 \right) \end{aligned}$$

Noting that  $m_d = \rho_l V = \rho_l \frac{4}{3} \pi \left(\frac{D}{2}\right)^3$ , then

$$\dot{m}_d = \rho_l 4\pi \frac{D^2}{8} \frac{dD}{dt}$$

Keeping in mind that  $\dot{m}_d = -\dot{m}$  because the mass of the droplet decreases with each layer of mass on the surface moving outward and away as a vapor,

$$\begin{aligned} \frac{dD^2}{dt} &= -\frac{8K_g}{c_{pg}\rho_l} \ln \left( 1 + \frac{c_{pg}}{h_{fg}} (T_\infty - T_b) \right) \\ \frac{dD^2}{dt} &= -\frac{8K_g}{\rho_l c_{pg}} \ln \left[ 1 + \frac{c_{pg}(T_\infty - T_b)}{h_{fg}} \right] = -\text{const.} \equiv -k_{\text{evap}} \end{aligned}$$

The decrease in  $D^2$  is constant! Then simply  $D^2(t) = D_0^2 - k_{\text{evap}} t$ . Thus, the lifetime of the droplet  $t_d$  is

$$t_d = \frac{D_0^2}{k_{\text{evap}}} = \frac{D_0^2}{\frac{8k_g}{\rho_l c_{pg}} \ln \left[ 1 + \frac{c_{pg}(T_\infty - T_b)}{h_{fg}} \right]}$$

21.1.2. *Gas Phase Composition.*

$$(82) \quad \boxed{\frac{d\phi_g}{dx} = \frac{1}{(F/O)_{\phi=1}} \frac{1}{\dot{m}_{\text{Ox}}(0)} \frac{d\dot{m}_g}{dx}}$$

22. DROPLET MODEL; BURNING DROPLETS

cf. 20160129 Dr. Jay Polk Ae121b Winter 2016

What we want to solve for:

$$\begin{aligned} \dot{m}_f &\implies \text{fuel production rate} \\ D(t) &\implies \text{applet lifetime} \end{aligned}$$

5 unknowns:

$$\begin{aligned} T_f \\ T_s \\ Y_{f,s} \\ r_f \\ \dot{m}_f \end{aligned}$$

The assumptions for this model is also given on pp. 379 of Turns (2011) [13], with the section **Simple Model of Droplet Burning**.

Note Assumption 3: “The fuel is a single-component liquid with zero solubility for gases.” Thus, there cannot be any net mass flux or mass flow of products inward from the flame to the droplet surface, since products cannot dissolve in the liquid by assumption 3.

Thus, in the inner zone  $r \in [r_s, r_f]$ ,  $\dot{m}_{\text{Products}} \equiv \dot{m}_{\text{Pr}}$ . However, there is products in the inner zone, that form a stagnant film through which the fuel vapor flows through.

**22.1. Species Conservation (Inner Zone).**  $F$  denotes fuel. We're thinking about the diffusion of fuel vapor, according to Fick's law, and so  $\mathcal{D}$  is the diffusivity constant in Fick's law, and is for binary diffusion (fuel vapor and combustion products), and for fuel vapor, in this case.  $\rho$  is for the fuel density.

Let  $Z_F := \frac{1}{4\pi\rho\mathcal{D}}$ .

The boundary conditions (BCs) are

$$\begin{aligned} Y_F(r_s) &= Y_{F,s}(T_s) \\ Y_F(r_f) &= 0 \\ Y_F(r) &= 1 - \frac{(1 - Y_{F,s}) \exp(-Z_F \dot{m}_F / r)}{\exp(-Z_F \dot{m}_F / r_s)} \end{aligned}$$

Applying boundary condition  $Y_F(r_f) = 0$ ,

$$\implies \frac{(1 - Y_{F,s}) \exp(-Z_F \dot{m}_F / r_f)}{\exp(-Z_F \dot{m}_F / r_s)} = 1$$

So

$$(83) \quad \boxed{Y_{F,s}(r_s) = 1 - \frac{\exp(-Z_F \dot{m}_F / r_s)}{\exp(-Z_F \dot{m}_F / r_f)}}$$

**22.2. Species Conservation (Outer Zone).** Assume stoichiometric combustion.

1 kg fuel +  $\nu$  kg oxidizer  $\rightarrow$   $\underbrace{(\nu + 1)}_{\text{get from stoichiometry}}$  kg of products.

EY : 20160216 Note that by mass conservation,

$$\begin{aligned} m_f + m_{\text{Ox}} &= m_p \implies 1 + \frac{m_{\text{Ox}}}{m_f} = \frac{m_p}{m_f} \\ \nu_f M_f + \nu_{\text{Ox}} M_{\text{Ox}} &= m_p \implies 1 + \frac{\nu_{\text{Ox}}}{\nu_f} \frac{M_{\text{Ox}}}{M_f} = \frac{m_p}{\nu_f M_f} \end{aligned}$$

The boundary conditions (BCs) are

$$\begin{aligned} Y_{\text{Ox}}(r_f) &= 0 \\ Y_{\text{Ox}}(\infty) &= 1 \end{aligned}$$

With

$$Y_{\text{Ox}}(r) = \nu \left[ \frac{\exp(-Z_F \dot{m}_F / r)}{\exp(-Z_F \dot{m}_F / r_f)} - 1 \right]$$

then, from  $Y_{\text{Ox}}(\infty) = 1$ ,

$$(84) \quad \boxed{\exp(Z_F \dot{m}_F / r_f) = \frac{\nu + 1}{\nu}}$$

One needs to take care in using the above. For instance, for the *outer zone*, by Assumption 5, *only the oxidizer and products* are there. So in

$$\dot{m}_F = 4\pi r^2 \frac{\rho_{\mathcal{D}}}{\nu + Y_{\text{Ox}}} \frac{dY_{\text{Ox}}}{dr}$$

and in solving it

$$\begin{aligned} \implies \frac{dr}{r^2} &= \frac{4\pi\rho\mathcal{D}\dot{m}_F}{\nu + Y_{\text{Ox}}} dY_{\text{Ox}} \\ \implies \frac{-1}{r} + \frac{1}{r_f} &= \frac{4\pi\rho\mathcal{D}}{\dot{m}_F} (\ln(\nu + Y_{\text{Ox}}) - \ln(\nu + Y_{\text{Ox}}(r_f))) \end{aligned}$$

then  $Z_f = \frac{1}{4\pi\rho\mathcal{D}}$  was used, but in this case, with binary diffusion between species  $A$ , which is the oxidizer, Ox, in this case, and species  $B$  which is the products, Pr, in this case, then  $\rho = \rho_A + \rho_B$  is for the total density of Ox and Pr in this case, not fuel F and Pr. Also,  $\mathcal{D}$  in this case refers to the binary diffusion of Ox and Pr, not to F and Pr, as in the inner zone, previously.

**22.3. Energy conservation at droplet.** Chemical reactions confined to occur only at the boundary i.e. the flame sheet (flame front). The reaction rate term is 0 both inside and outside the flame.

Thus, the energy equation for droplet evaporation is used.

$Z_T = \frac{c_{pg}}{4\pi k_g} = Z_F$  if Le = 1, Lewis number is 1.

$$(85) \quad \underbrace{\dot{q}_{i-l}}_{\text{interface into liquid}} \quad \underbrace{\dot{q}_{g-i}}_{\text{gas to interface}} = - \left[ -k_g 4\pi r^2 \frac{dT}{dr} \Big|_{r_s} \right]$$

$$\boxed{\frac{c_{pg}(T_f - T_s)}{(q_{i-l} + h_{fg})} \frac{\exp(-Z_T \dot{m}_F / r_s)}{[\exp(-Z_T \dot{m}_F / r_s) - \exp(-Z_T \dot{m}_F / r_f)]} + 1 = 0}$$

**22.4. Energy Conservation at Flame Front.** Chemical energy released at flame taken account by absolute enthalpy fluxes for fuel, oxidizer, and products:

$$\dot{m}_F h_F + \dot{m}_{\text{Ox}} h_{\text{Ox}} - \dot{m}_{\text{Pr}} h_{\text{Pr}} = \dot{Q}_{f-i} + \dot{Q}_{f-\infty}$$

Enthalpies are defined as

$$\begin{aligned} h_F &:= h_{f,F}^\circ + c_{pg}(T - T_{\text{ref}}) \\ h_{\text{Ox}} &:= h_{f,\text{Ox}}^\circ + c_{pg}(T - T_{\text{ref}}) \\ h_{\text{Pr}} &:= h_{f,\text{Pr}}^\circ + c_{pg}(T - T_{\text{ref}}) \end{aligned}$$

Heat of combustion  $\Delta h_c$  per unit mass of fuel given by

$$\Delta h_c(T_{\text{ref}}) = (1)h_{f,F}^\circ + (\nu)h_{f,\text{Ox}}^\circ - (1 + \nu)h_{f,\text{Pr}}^\circ$$

Note that although products exist in inner region,  $\nexists$  net flow of products between droplet surface and flame; thus all products flow radially outward away from flame, i.e.

$$\begin{array}{ccc} r_s & \xrightarrow{\dot{m}_F} & r_f \\ & & \nwarrow \dot{m}_{\text{Pr}} \\ & & \text{Thus } \dot{m}_F + \dot{m}_{\text{Ox}} = -\dot{m}_{\text{Pr}}, \text{ or } \dot{m}_F + \nu \dot{m}_F = -\dot{m}_{\text{Pr}} \end{array}$$

or thus  $\dot{m}_F(1 + \nu) = -\dot{m}_{\text{Pr}}$

Hence

$$\begin{aligned} \dot{m}_F [h_F + \nu h_{\text{Ox}} - (\nu + 1)h_{\text{Pr}}] &= \dot{Q}_{f-i} + \dot{Q}_{f-\infty} \\ \implies \dot{m}_F \Delta h_c + \dot{m}_F c_{pg} [(T_f - T_{\text{ref}}) + \nu(T_f - T_{\text{ref}}) - (\nu + 1)(T_f - T_{\text{ref}})] &= \dot{Q}_{f-i} + \dot{Q}_{f-\infty} \end{aligned}$$

Since we assume  $c_{pg}$  constant, then  $\Delta h_c$  independent of temperature; thus choose flame temperature as reference state,

$$\implies \dot{m}_F \Delta h_c = \dot{Q}_{f-i} + \dot{Q}_{f-\infty}$$

Therefore, the trick that we can employ is to simply only consider the  $\Delta h_c$  at the flame temperature.

$$\begin{array}{ccc} \text{interface} & \xleftarrow{\dot{q}_{f-i}} & \parallel \text{flame} \parallel \xrightarrow{\dot{q}_{f-\infty}} \infty \\ \text{i} & \xrightarrow{\dot{m}_f h_f} & \parallel f \parallel \xleftarrow{\dot{m}_{\text{Ox}} h_{\text{Ox}} = \nu \dot{m}_f h_{\text{Ox}}} \infty \\ & & \parallel f \parallel \xrightarrow{\dot{m}_p h_p = (\nu + 1) \dot{m}_f h_p} \infty \end{array}$$

Thus

$$\begin{aligned}\dot{m}_f[h_f + \nu h_{\text{Ox}} - (\nu + 1)h_f] &= \dot{q}_{f-i} + \dot{q}_{f-\infty} \\ \Delta h_c(T_{\text{ref}}) &= h_{f,f}^\circ + \nu h_{\text{Ox},f}^\circ - (\nu + 1)h_{f,f}^\circ \\ \dot{m}_f \Delta h_c &= \dot{q}_{f-i} + \dot{q}_{f-\infty}\end{aligned}$$

(92)

Using  $\left.\frac{dT}{dr}\right|_{\text{rs}}$  from inner+outer energy conservation.

(86)

$$\frac{c_{pg}}{\Delta h_c} \left[ \frac{(T_s - T_f) \exp(-Z_T \dot{m}_F / r_f)}{\exp(-Z_T \dot{m}_F / r_s) - \exp(-Z_T \dot{m}_F / r_f)} - \frac{(T_\infty - T_f) \exp(-Z_T \dot{m}_F / r_f)}{[1 - \exp(-Z_T \dot{m}_F / r_f)]} \right] - 1 = 0$$

22.5. **Liquid-Vapor Equilibrium at Droplet Surface.** Clausius-Clapeyron Eqn.

$$p_v = p_0 \exp \left[ \frac{h_{fg}}{R} \left( \frac{1}{T_0} - \frac{1}{T} \right) \right]$$

can rewrite as

$$p_v = A \exp \left( \frac{-B}{T_s} \right)$$

$A, B$  are constants for a given liquid.

So fuel partial pressure at surface  $\simeq$  equilibrium vapor.

pressure  $p_{F,s} = A \exp \left( \frac{-B}{T_s} \right)$

Fuel mole fraction  $X_{F,s} = \frac{p_{F,s}}{p}$  and

(87)

$$Y_{F,s} = \frac{N_F m_F}{N_F m_F + N_P m_P} = X_F \frac{m_F}{X_F m_F + (1 - X_F) m_P}$$

$$Y_{F,s} = \frac{A \exp(-B/T_s) m_F}{A \exp(-B/T_s) m_F + [p - A \exp(-B/T_s)] m_P}$$

22.5.1. *Empirical relations.* Turns (2011) [13] quotes Law and Williams (1972) (cf. Law, C.K., and Williams, F.A., “Kinetics and Convection in the Combustion of Alkane Droplets,” *Combustion and Flame*, 19(3): 393-406 (1972)) for empirical relations to use for  $k_g$  and  $c_{pg}$  (EY : 20160227 I need to check this article out myself).

Therefore,

(88)

$$\begin{aligned}c_{pg} &= c_{pF}(\bar{T}) \\ k_g &= 0.4k_F(\bar{T}) + 0.6k_{\text{Ox}}(\bar{T}) \\ \rho_l &= \rho_l(T_s)\end{aligned}$$

Also, from 20160129 Dr. Polk, Ae121b Winter 2016,  $\bar{T} = (T_s + T_f)/2$ , and a good initial guess for  $T_s$  and  $T_f$  are  $T_b(P)$  and  $T_{\text{ad}}$ , the adiabatic flame temperature for stoichiometric mixture.

After algebraic manipulations (which can be done in `BurningDroplet.py` with `sympy` in Python, instead of by hand), then we have 5 equations

(89)

$$\dot{m}_F = \frac{4\pi k_g r_s}{c_{pg}} \ln [1 + B_{oq}]$$

(90)

$$T_f = \frac{q_{i-l} + h_{fg}}{c_{pg}(1 + \nu)} [\nu B_{oq} - 1] + T_s$$

(91)

$$r_f = r_s \frac{\ln(1 + B_{oq})}{\ln[(\nu + 1)/\nu]}$$

•

$$Y_{F,s} = \frac{B_{oq} - 1/\nu}{B_{oq} + 1}$$

•

(93)

$$T_s = \frac{-B}{\ln \left[ \frac{-Y_{F,s} p m_P}{A(Y_{F,s} m_F - Y_{F,s} m_P - m_F)} \right]}$$

with

$$B_{oq} := \frac{\Delta h_c / \nu + C_p (T_\infty - T_s)}{q_{i-l} + h_{fg}}$$

Keep in mind the constant decrease in the size of the droplet, measured by diameter squared, or  $D^2$ ,

$$\frac{dD^2}{dt} = -\kappa$$

with

$$\kappa := \frac{8k_g}{\rho_l c_{pg}} \ln(1 + B_{oq})$$

I will compare with Turns (2011) [13], from pp. 378, Section **Simple Model of Droplet Burning** of Chapter 10 Droplet Evaporation and Burning.

Turns (2011) [13] lists 10 assumptions on pp. 379 for the Simple Model of Droplet Burning. Beginning with his fifth assumption,

The gas phase consists of only 3 species: fuel vapor, oxidizer, and combustion products. The gas phase region is divided into 2 zones, inner zone, and outer zone. So for this spherically symmetric problem, parametrize space by radius  $r \in \mathbb{R}^+$  from the center of the (fuel) droplet.

Assumption 3 of Turns (2011) [13]: Fuel is single-component liquid with 0 solubility for gases. Phase equilibrium at liquid-vapor interface (which I believe is at the surface of the droplet,  $r_s$ ).

Assumption 10 of Turns (2011) [13]: “Liquid fuel droplet is the only condensed phase; no soot or liquid water is present.”

$Y_F \equiv$  fuel vapor mass fraction.

22.5.2. *Inner zone.* Define the inner zone  $r \in [r_s, r_f]$  between droplet surface  $r_s$  and flame front  $r_f$ .

Inner zone only contains fuel vapor and combustion products: binary diffusion prevails.

$T_s \equiv$  droplet surface temperature.

$T_f \equiv$  flame temperature.

$Y_{F,s} = Y_F(r = r_s) \equiv Y_F(r_s) =$  fuel vapor mass fraction at the droplet surface

$Y_F(r_f) = 0$  means that all the fuel is consumed by the time we reach the flame front.

Turns mentions that “a more elegant approach” is described in Kuo’s **Principles of Combustion**, that “combines the species and energy equations to create a conserved-scalar variable.

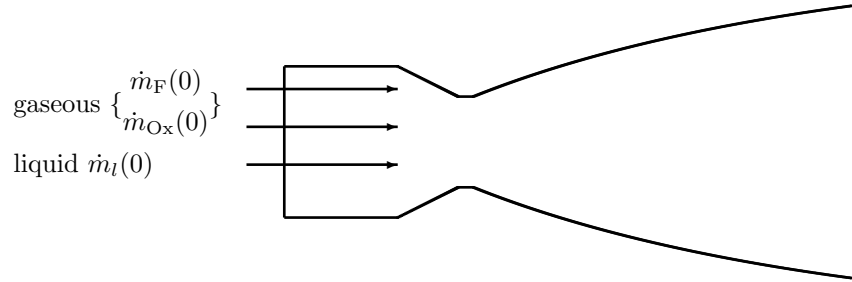
## 23. COMBUSTION CHAMBER FLOW MODEL

cf. 20160202 Ae121b Dr. Polk.

Parametrize the axisymmetric axis  $x \in \mathbb{R}$ , with  $x = 0$  being the inlet exit into the combustion chamber.

Consider the (mass) inlet flow of the oxidizer,  $\dot{m}_{\text{Ox}}(x = 0) \equiv \dot{m}_{\text{Ox}}(0)$  and the (mass) inlet flow of the fuel  $F$ ,  $\dot{m}_F(x = 0) \equiv \dot{m}_F(0)$ . Here, fuel is a vapor (i.e. gas).

Suppose some of the fuel from the inlet will come in as a liquid  $l$ ; denote the inlet flow of liquid fuel as  $\dot{m}_l(0)$ .



Consider the following assumptions:

- There are only **2 phases**: the gaseous phase of fuel + oxidizer, and the liquid phase of liquid fuel.
- Assume 1-dimensional flow, with no diffusion in the bulk fluid.
- Assume constant pressure in the combustion chamber, and adiabatic flow, with no work on the flow.
- Fuel is injected as a *monodisperse* spray, meaning that all the droplets have the same initial diameter.
- Gas phase is in equilibrium at even the axial position  $x$ .

The parameters of interest along the axis  $x$  are the following:

$$\begin{array}{ll} \dot{m}_g(x) & \dot{m}_l(x) \\ T_g(x) & D(x) \\ \phi_g(x) & v_d(x) \\ v_g(x) & \end{array}$$

**23.1. Inlet conditions.**  $\dot{m}_g(0)$  refers to the inlet.

$\dot{m}_f(0)$  refers to flow rate of fuel injected as gas.

Then by accounting for the total mass of the gas injected at the inlet, being part gaseous fuel,  $f$ , and gaseous oxidizer, Ox,

$$\dot{m}_g(0) = \dot{m}_f(0) + \dot{m}_{\text{Ox}}(0)$$

By definition,

$$\left(\frac{F}{O}\right)_x = \frac{\dot{m}_F(x)}{\dot{m}_{\text{Ox}}(x)}$$

Now

$$\phi_g(x = 0) = \phi_g(0) = \frac{\dot{m}_f(0)/\dot{m}_{\text{Ox}}(0)}{(F/O)_{\Phi=1}} = 0.45$$

by definition of an *equivalence ratio*,  $\phi_g$  for the gas  $g$ , and so we can determine  $\dot{m}_f(0)$  after doing some algebra

$$(94) \quad \dot{m}_f(0) = \dot{m}_{\text{Ox}}(0) \left(\frac{F}{O}\right)_{\Phi=1} \phi_g(0)$$

Now, if given an overall equivalence ratio,  $\phi_{\text{overall}}$  at a point  $x$  along the axis,

$$(95) \quad \phi_{\text{overall}}(x = 0) = \frac{\dot{m}_{\text{fuel}}(0)/\dot{m}_{\text{Ox}}(0)}{(F/O)_{\Phi=1}} = \frac{\frac{\dot{m}_f(0)}{\dot{m}_{\text{Ox}}(0)} + \frac{\dot{m}_l(0)}{\dot{m}_{\text{Ox}}(0)}}{(F/O)_{\Phi=1}} = \phi_g(0) + \frac{\frac{\dot{m}_l(0)}{\dot{m}_{\text{Ox}}(0)}}{(F/O)_{\Phi=1}}$$

Clearly, the total mass of fuel inputted in is the mass of the fuel as a gas and mass of the fuel as a liquid (droplets):

$$(96) \quad \dot{m}_{\text{Fuel}}(0) = \dot{m}_f(0) + \dot{m}_l(0)$$

By the physical parameters of a given setup, such as injection velocity of the gas  $v_g(x = 0) \equiv v_g(0)$  and total fuel injector cross-sectional area,  $A_{\text{inlet Fuel}}$

$$(97) \quad \dot{m}_g(0) = \rho_g v_f(0) A_{\text{inlet Fuel}} = \rho_g (T_{\text{inlet}}, P) v_f(0) A_{\text{inlet Fuel}}$$

where we remember that  $\rho_f$  is a function of inlet temperature  $T_{\text{inlet}}$  and pressure  $P$  of the combustion chamber.  $\rho_g$  can be obtained by setting the state of the gas with *Cantera*, with the appropriate equivalence ratio for the gases  $\phi_g(0)$ .

From Eq. [94](#), we can obtain  $\dot{m}_{\text{Ox}}(0)$  with the help of Eq. [97](#), namely

$$\dot{m}_g(0) = \dot{m}_{\text{Ox}}(0) \left(1 + \left(\frac{F}{O}\right)_{\Phi=1} \phi_g(0)\right)$$

And then, from the first equality of Eq. [95](#), and Eq. [96](#), then  $\dot{m}_l(0)$  can be obtained, namely

$$\phi_{\text{overall}}(0)(F/O)_{\Phi=1} = \frac{\dot{m}_f(0) + \dot{m}_l(0)}{\dot{m}_{\text{Ox}}(0)}$$

**23.2. Gas phase continuity.** Since for a gas,  $pV = N\tau$ , so

$$p = \frac{N}{V}\tau = \frac{MN}{V} \frac{1}{M} \frac{T}{k_B} = \rho \frac{1}{M k_B} T = \rho R_M T$$

$$(98) \quad \boxed{v_g = \frac{\dot{m}_g}{\rho_g A} = \frac{\dot{m}_g R_U T_g}{m_g P A}}$$

$$\text{i.e. } v_g = \frac{\dot{m}_g \tau_g}{M_g p A}$$

**23.3. Gas phase Energy Conservation.** For  $H_{\text{tot}} = H_{\text{tot}}(\tau, p, \{N_i\}, x) \equiv H_{\text{tot}}(x) \in C^\infty(\Sigma \times N) = C^\infty(\Sigma \times \mathbb{R})$ ,  $H_{\text{tot}} = m_g h_g + m_l h_l$ .

For  $\dot{H}_{\text{tot}} = \dot{m}_g h_g + \dot{m}_l h_l$  where I take into consideration the steady state assumption, then by

$$Q = 0, dH = 0, \text{ for } dp = 0 \implies \frac{d}{dx}(\dot{m}_g h_g) + \frac{d}{dx}(\dot{m}_l h_l) = 0$$

Assume constant  $T$  droplets, so  $h_l = \text{constant}$

$$(99) \quad \frac{dh_g}{dx} = \frac{-1}{\dot{m}_g} \left[ h_g \frac{d\dot{m}_g}{dx} + h_l \frac{d\dot{m}_l}{dx} \right]$$

Since  $h_g = h_g(T_g, p, \phi_g) \in C^\infty(\Sigma)$ ,

$$(100) \quad \frac{dh_g}{dx} = \frac{\partial h_g}{\partial T} \frac{dT}{dx} + \frac{\partial h_g}{\partial \phi_g} \frac{d\phi_g}{dx} \quad (p \text{ is const.})$$

Equating Eqns. [99](#), [100](#), then

$$\frac{dT_g}{dx} = \left[ \frac{-1}{\dot{m}_g} \left( h_g \frac{d\dot{m}_g}{dx} + h_l \frac{d\dot{m}_l}{dx} \right) - \frac{\partial h_g}{\partial \phi_g} \frac{d\phi_g}{dx} \right] / \frac{\partial h_g}{\partial T_g}$$

Using mass conservation on fuel,  $\frac{d\dot{m}_g}{dx} = -\frac{d\dot{m}_l}{dx}$ , so

$$(101) \quad \boxed{\frac{dT_g}{dx} = \left[ \frac{-1}{\dot{m}_g} \frac{d\dot{m}_g}{dx} (h_g - h_l) - \frac{\partial h_g}{\partial \phi_g} \frac{d\phi_g}{dx} \right] / \frac{\partial h_g}{\partial T_g}}$$



### 23.4. Droplet momentum conservation.

$$F_d = m_d \frac{dv_d}{dt} = m_d \frac{dx}{dt} \frac{dv_d}{dx} = m_d v_d \frac{dv_d}{dx}$$

$$F_d \text{ for a sphere} = C_d \rho_g \frac{v_{\text{rel}}^2}{2} \frac{\pi D^2}{4}$$

Recall the Reynolds number for a sphere moving through a fluid,

$$\text{Re} = \frac{D_{\text{droplet}} u_{\text{rel}} \rho_{\text{fluid}}}{\mu_{\text{fluid}}}$$

where  $\mu$  viscosity of the fluid,  $\rho_{\text{fluid}} = \rho_g$ ,  $\mu_{\text{fluid}} = \mu_g$ . Keep in mind that  $[D_{\text{droplet}} u_{\text{rel}} \rho_{\text{fluid}}] = m \cdot \frac{m}{s} \cdot \frac{\text{kg}}{m^3} = \frac{kg}{m \cdot s}$ .

$$(102) \quad \begin{aligned} \frac{dD^2}{dx} &= \frac{-K}{v_d} \\ \frac{dv_d}{dx} &= \frac{3C_D \rho_g (v_g - v_d) |v_g - v_d|}{4\rho_l D v_d} \\ \frac{dT_g}{dx} &= \left[ \frac{-1}{\dot{m}_g} \frac{d\dot{m}_g}{dx} (h_g - h_l) - \frac{\partial h_g}{\partial \phi_g} \frac{d\phi_g}{dx} \right] / \frac{\partial h_g}{\partial T_g} \end{aligned}$$

droplet  $\in \Sigma \times \Gamma(TN)$

$(\rho_l, p_l) \times (x, v_d(x))$

assume constant  $T$  droplets.

So I claim by clausius Clapeyron  $T_l = T_{\text{boil}}$

Let's take a look at Eq. [102](#).

- Looking at  $\frac{dD^2}{dx} = \frac{-K}{v_d}$ , and looking at what each of the factor's respective formulae,

$$\frac{dD^2}{dx} = \frac{-K}{v_d} \implies \begin{aligned} K &= \frac{8k_g}{\rho_l C_{pg}} \ln(1 + B_{oq}) \\ B_{oq} &= \frac{\Delta h_c / \nu + C_{pg}(T_\infty - T_s)}{q_{i-l} + h_{fg}} \end{aligned}$$

Then, examining the terms *related to gases* that depend on  $\Sigma$  and spatial manifold  $N$ , parametrize by  $x \in \mathbb{R}$ ,

$$\begin{aligned} k_g &= k_g(T_g, P, \phi_g) \in C^\infty(\Sigma) \\ k_g(T_g(x), P, \phi_g(x)) &\in C^\infty(\Sigma_x) \implies k_g(x) \in C^\infty(N) \end{aligned}$$

with  $\Sigma \rightarrow N$  being a fibered bundle, with fibers  $\Sigma_x$ .

Likewise,

$$\begin{aligned} C_{Pg} &= C_{Pg}(T_g, P, \phi_g) \in C^\infty(\Sigma) \\ C_{Pg}(T_g(x), P, \phi_g(x)) &\in C^\infty(\Sigma_x) \text{ and } C_{Pg}(x) \in C^\infty(N) \end{aligned}$$

In  $B_{oq}$ , for term  $h_{fg}$ , for the heat of formation,

$$\begin{aligned} h_{fg} &= h_{fg}(T^0, P, \phi_g) \in C^\infty(\Sigma) \\ h_{fg}(T^0, P, \phi_g(x)) &\in C^\infty(\Sigma_x) \text{ and } h_{fg}(x) \in C^\infty(N) \end{aligned}$$

and since assuming  $C_{PF} = C_{POx} = C_{pg}$  then

$$\Delta h_c = \Delta h_c(T) \in C^\infty(\Sigma) \text{ and } \Delta h_c(T = T_{\text{ref}}) \equiv \Delta h_c(T_{\text{ref}}) = \Delta h_c(T)$$

with the upshot that the “heat of combustion” is the same at  $T = T_{\text{ref}}$  and at arbitrary  $T$ .

For terms involving the droplet, the “*liquid*”,

$$\rho_l = \rho_l(T) \in C^\infty(\Sigma) \rho_l(T(x)) \in C^\infty(\Sigma_x) \text{ and } \rho_l(x) \in C^\infty(N)$$

In our particular case, we assume the droplet has uniform temperature, which happens to be the temperature at the surface, which then happens to be dependent on pressure, by the Clausius-Clapeyron relation. Keep in mind that the relation  $\rho_l = \rho_l(T)$  is from the relation involving the volumnic thermal expansion  $\alpha$  of a *liquid*:  $\rho = \frac{\rho_{T_0}}{1 + \alpha(T - T_0)}$ , which, I think, is an empirical relation, as  $\alpha$  needs to be measured and is a parameter to input in.

$$\rho_l = \rho_l(T_s) = \rho_l(T_b) \in C^\infty(\Sigma)$$

- Looking at  $\frac{dv_d}{dx} = \frac{3C_D \rho_g (v_g - v_d) |v_g - v_d|}{4\rho_l D v_d}$ , and taking a look at terms involving the *gaseous* phase,

$$v_g \in TN(x, v_g) = v_g(x) \in T_x N$$

and

$$\begin{aligned} \rho_g &= \rho_g(T_g, P, \phi_g) \in C^\infty(\Sigma) \\ \rho_g(T_g(x), P, \phi_g(x)) &\in C^\infty(\Sigma_x) \text{ and } \rho_g(x) \in C^\infty(N) \end{aligned}$$

Looking at terms involving the *liquid* phase, the **droplet**,

$$\begin{aligned} v_d &\in TN \\ (x, v_d) &= v_d(x) \in T_x N \\ D &= D(x) \in C^\infty(N) \end{aligned}$$

Then, examining the formulae for the Reynolds number  $\text{Re}$  and  $C_D$ , affecting the drag on the droplet,

$$\text{Re}_{D, \text{rel}} = \frac{D |v_g - v_d| \rho_g}{\mu}$$

with viscosity  $\mu$  being a transport property of the *gas* that the droplet is moving around in, and

$$\begin{aligned} \mu &= \mu(T_g, P, \phi_g) \in C^\infty(\Sigma) \\ \mu(T_g(x), P, \phi_g(x)) &\in C^\infty(\Sigma_x) \text{ and } \mu(x) \in C^\infty(N) \end{aligned}$$

$C_D$  drag coefficient is fairly straightforward, depending upon  $x \in N$ , but through  $\text{Re}_{D, \text{rel}}$ , *only*:

$$C_D \simeq \frac{24}{\text{Re}_{D, \text{rel}}} + \frac{6}{1 + \sqrt{\text{Re}_{D, \text{rel}}}} + 0.4$$

- Looking at  $\frac{dT_g}{dx} = \left[ \frac{-1}{\dot{m}_g} \frac{d\dot{m}_g}{dx} (h_g - h_l) - \frac{\partial h_g}{\partial \phi_g} \frac{d\phi_g}{dx} \right] / \frac{\partial h_g}{\partial T_g}$ , the terms involving the *gas* phase are the following:

$$\begin{aligned} T_g &\in \Sigma \text{ and } T_g(x) \in \Sigma(x) \equiv \Sigma_x \text{ and} \\ T_g(x) &\in C^\infty(N) \\ \phi_g &= \phi_g(x) \in C^\sigma(\Sigma_x) \text{ and} \\ \phi_g(x) &\in C^\infty(N) \\ h_g &= h_g(T_g, P, \phi_g) \in C^\infty(\Sigma) \\ h_g(T_g(x), P, \phi_g(x)) &\in C^\infty(\Sigma_x) \text{ and } h_g(x) \in C^\infty(N) \end{aligned}$$

and the (physical) quantity that doesn't depend upon  $\Sigma$  for the *gaseous* phase is

$$\dot{m}_g = \dot{m}_g(x) \in C^\infty(N)$$

For the terms involving the droplet, or the *liquid* phase,

$$\begin{aligned} h_l &= h_l(T_l, P) \in C^\infty(\Sigma) \\ h_l(T_s, P) &= h_l(T_{\text{boil}}, P) \in C^\infty(\Sigma) \end{aligned}$$

23.4.1. *Considerations of the physical parameters involved.* Start with the “mass flow” rates that are smooth functions of space and  $N$ . There are 2 reactant *species* involved, the fuel and oxidizer, Ox, amongst 2 phases, gas and liquid, considered.

$$\dot{m}_g = \dot{m}_g(x), \dot{m}_f = \dot{m}_f(x), \dot{m}_{\text{Ox}} = \dot{m}_{\text{Ox}}(x), \dot{m}_l = \dot{m}_l(x) \in C^\infty(N)$$

Clearly, we have

$$(103) \quad \dot{m}_g = \dot{m}_f + \dot{m}_{\text{Ox}}$$

as gaseous oxidizer is injected in, and gaseous fuel is from the evaporation off the liquid droplet.

Modelling the liquid fuel being injected into the combustion chamber through an injection plate of total fuel injection cross-section  $A_{\text{tot fuel inj}}$ , at injection speed  $v_d(0)$ , then clearly

$$(104) \quad \dot{m}_l(x=0) = \dot{m}_l(0) = \rho_l v_d(0) A_{\text{tot fuel inj}} = \rho_l(T_b, P) v_d(0) A_{\text{tot fuel inj}}$$

I would argue that for the liquid density  $\rho_l = \rho_l(T_b, P)$ , which is completely specified by 2 thermodynamic quantities in  $\Sigma^l$ ,  $(T, P) \in \Sigma^l$ , the temperature is  $T_b$ , the boiling temperature specified by the Clausius-Clapeyron relation for a given combustion chamber pressure  $P$ . The inlet temperature is usually much higher than the boiling temperature of our fuel: we wouldn’t have any liquid, or droplets, if its local temperature is this inlet temperature. I will also argue that the liquid droplet is entirely in thermal equilibrium and at its surface, fuel molecules are evaporating away at  $T_b$ : so the entire liquid is at  $T_b = T_b(P)$ . This remains the case as the liquid droplet travels throughout the combustion chamber.

Out of this injection plate, for a droplet of initial size  $D_0$ , then for number of droplets per unit time emerging out of the injection plate,  $\dot{N}$ , clearly

$$\dot{m}_l(0) = \frac{\dot{N} \rho_l \pi D_0^3}{6}$$

Consider the so-called equivalence ratio; in this case the overall equivalence ratio  $\phi_{\text{overall}}$ . It is an interesting quantity that relates purely physical quantities, masses or number of particles (i.e. moles), to thermodynamic variables  $\{N_i\}_i \in \Sigma$ , i.e.  $\phi_{\text{overall}} \in C^\infty(\Sigma)$ . By definition of  $\phi_{\text{overall}}$ ,

$$(105) \quad \phi_{\text{overall}} := \frac{\frac{\dot{m}_f + \dot{m}_l}{\dot{m}_{\text{Ox}}}}{\left(\frac{F}{O}\right)_{\Phi=1}} = \frac{\frac{\dot{m}_g - \dot{m}_{\text{Ox}} + \dot{m}_l}{\dot{m}_{\text{Ox}}}}{\left(\frac{F}{O}\right)_{\Phi=1}} \implies \dot{m}_g = \left( \phi_{\text{overall}} \left( \frac{F}{O} \right)_{\Phi=1} + 1 \right) \dot{m}_{\text{Ox}} - \dot{m}_l$$

Also, by definition of the equivalence ratio of gaseous fuel to oxidizer (which is assumed to have completely vaporized into gas),  $\phi_g \in C^\infty(\Sigma)$ ,

$$(106) \quad \phi_g := \frac{\frac{\dot{m}_f}{\dot{m}_{\text{Ox}}}}{\left(\frac{F}{O}\right)_{\Phi=1}} = \frac{\frac{\dot{m}_g - \dot{m}_{\text{Ox}}}{\dot{m}_{\text{Ox}}}}{\left(\frac{F}{O}\right)_{\Phi=1}} \implies \dot{m}_g = \left( \left( \frac{F}{O} \right)_{\Phi=1} \phi_g + 1 \right) \dot{m}_{\text{Ox}}$$

By equating the expressions for  $\dot{m}_g$  in Eqns. 105, 106, then

$$\dot{m}_{\text{Ox}} = \frac{\dot{m}_l}{(\phi_{\text{overall}} - \phi_g) \left( \frac{F}{O} \right)_{\Phi=1}}$$

Notice that  $\phi_{\text{overall}} > \phi_g \geq 0$ .

I haven’t mentioned where this is all occurring. Here is one important fact about the physical setup:

$$\dot{m}_{\text{Ox}}(x) = \dot{m}_{\text{Ox}}(0) \quad \forall x \in N$$

i.e. the oxidizer is constantly being injected in and flowing uniformly through the chamber.

At  $x = 0$ , the start of the combustion chamber,

$$(107) \quad \dot{m}_{\text{Ox}}(0) = \frac{\dot{m}_l(0)}{(\phi_{\text{overall}}(x=0) - \phi_g(x=0)) \left( \frac{F}{O} \right)_{\Phi=1}}$$

Let’s take stock of the quantities and its properties that we’ve discussed so far:

$$\begin{aligned} \dot{N}(x) &= \dot{N}(0) & \forall x \in N & \quad (\text{droplets don't split or breakup, and don't combine together}) \\ \dot{m}_{\text{Ox}}(x) &= \dot{m}_{\text{Ox}}(0) & \forall x \in N & \quad (\text{constant injection of oxidizer}) \end{aligned}$$

$$\dot{m}_g(x) = \left( \left( \frac{F}{O} \right)_{\Phi=1} \phi_g(x) + 1 \right) \dot{m}_{\text{Ox}}(0) = \dot{m}_g(\phi_g(x)) \in C^\infty(\Sigma_x)$$

where we notice that  $\dot{m}_g$  is a function of space  $x \in N$  or a function of thermodynamic property  $\phi_g(x)$ . Also note that the quantity  $\left( \frac{F}{O} \right)_{\Phi=1}$ , the so-called *stoichiometric* ratio between the *mass* of fuel to *mass* of oxidizer that combines together in a reaction that goes to 100% completion (hence “stoichiometric”); once the chemical reaction is chosen, this quantity remains *fixed* throughout  $x \in N$ .

Thus, so far, we should carry forward these 3 physical quantities, that are completely determined by inlet conditions:

$$(\dot{m}_l(0), \dot{N}(0), \dot{m}_{\text{Ox}}(0)) \in C^\infty(N_{x=0}) \times C^\infty(N_{x=0}) \times C^\infty(N_{x=0})$$

For  $x \in N$  (i.e. along the rest of the combustion chamber), and recalling this derivation,

$$\frac{d}{dx} \dot{m}_l = \frac{d}{dx} (\dot{N} \rho_l \frac{\pi D^3}{6}) = \frac{3}{2} \frac{\dot{m}_l(0)}{D_0^3} D \frac{dD^2}{dx} \xrightarrow{\frac{dD^2}{dx} = \frac{-K}{v_d}} \frac{-3}{2} \frac{\dot{m}_l(0)}{D_0^3} D \frac{K}{v_d}$$

*Droplet evaporation off surface* is expressed as:

$$\frac{d\dot{m}_g}{dx} = -\frac{d\dot{m}_g}{dx} \implies \frac{d\dot{m}_g}{dx} = \frac{3}{2} \frac{\dot{m}_l(0)}{D_0^3} D \frac{K}{v_d}$$

and thus, taking the spatial derivative of both sides of Eq. 106

$$\frac{d\phi_g}{dx} = \frac{1}{\left( \frac{F}{O} \right)_{\Phi=1} \dot{m}_{\text{Ox}}(0)} \frac{d\dot{m}_g}{dx} = \frac{1}{\left( \frac{F}{O} \right)_{\Phi=1} \dot{m}_{\text{Ox}}(0)} \frac{3}{2} \frac{\dot{m}_l(0)}{D_0^3} D \frac{K}{v_d}$$

By gas phase continuity,

$$v_g = v_g(x) = \frac{\dot{m}_g}{\rho_g A_{cc}} = \frac{\dot{m}_g(x) R_U T_g(x)}{MW_g P A_{cc}} = \frac{\left( \left( \frac{F}{O} \right)_{\Phi=1} \phi_g(x) + 1 \right) \dot{m}_{\text{Ox}}(0) R_U T_g(x)}{MW_g P A_{cc}}$$

We can obtain  $\dot{m}_l = \dot{m}_l(x)$  (i.e. as a function of  $x \in N$ ):

$$\dot{m}_l(x) = \frac{\dot{N} \rho_l(T_b, P) \pi D^3}{6}$$

Notice that it is a function of  $D = D(x) \in C^\infty(N)$ .

Now consider that

$$\frac{d\phi_{\text{overall}}}{dx} = \frac{\frac{d\dot{m}_g}{dx} - 0 + \frac{d\dot{m}_l}{dx}}{\dot{m}_{\text{Ox}} \left( \frac{F}{O} \right)_{\Phi=1}} = 0$$

i.e.  $\phi_{\text{overall}}$  is constant throughout the chamber! This could be attributed to overall mass conservation, but it’s explicitly shown above.

Therefore

$$(108) \quad \begin{aligned} \dot{m}_g(x) &= \left( \phi_{\text{overall}} \left( \frac{F}{O} \right)_{\Phi=1} + 1 \right) \dot{m}_{\text{Ox}}(0) - \dot{m}_l(x) \\ \implies \phi_g(x) &= \left( \frac{\dot{m}_g(x)}{\dot{m}_{\text{Ox}}(0)} - 1 \right) \frac{1}{\left( \frac{F}{O} \right)_{\Phi=1}} \end{aligned}$$

and so  $\phi_g(x)$  can be obtained from this formula above,  $\forall x \in N$ .

It may be interesting to notice that, and to take stock of what we’ve uncovered so far, that the dynamic quantities (dynamic in that they’ll be involved in our system of ordinary differential equations (ODEs))

$$(T_g(x), D(x), v_d(x)) \in C^\infty(\Sigma^g) \times C^\infty(N) \times TN$$

along with initial (inlet) conditions

$$(\dot{m}_{\text{Ox}}(0), \dot{N}, \dot{m}_l(0)) \in C^\infty(N_0) \times C^\infty(N_0) \times C^\infty(N_0)$$

completely determines, algebraically, the quantities

$$(\phi_g(x), \dot{m}_g(x), v_g(x)) \in C^\infty(\Sigma^g) \times C^\infty(N) \times TN, \quad \dot{m}_l(x) \in C^\infty(N)$$

all throughout  $x \in N$ .

Alternatively, and even *better*, keeping the assumption that  $\dot{\mathcal{N}}$  is constant, we should eliminate  $\dot{\mathcal{N}}$  algebraically and write only in terms of the flow rates: starting from  $\dot{m}_l(0) = \rho_l(T_{\text{boil}}, P)v_d(0)A_{\text{tot fuel inj}}$ , Eq. [104](#),

$$\dot{m}_l(0) = \frac{\dot{\mathcal{N}}\rho_l(T_{\text{boil}}, P)\pi D_0^3}{6} \text{ or } \frac{\dot{m}_l(0)}{D_0^3} = \frac{\dot{\mathcal{N}}\rho_l(T_{\text{boil}}, P)\pi}{6}$$

and thus

$$\dot{m}_l(x) = \dot{m}_l(0) \left( \frac{D}{D_0} \right)^3$$

From  $\dot{m}_{\text{Ox}}(0) = \frac{\dot{m}_l(0)}{(\phi_{\text{overall}}(0) - \phi_g(0))\left(\frac{F}{O}\right)_{\Phi=1}}$ , Eq. [107](#) So then

$$\dot{m}_g(0) = (\phi_{\text{overall}} \left( \frac{F}{O} \right)_{\Phi=1} + 1) \dot{m}_{\text{Ox}}(0) - \dot{m}_l(0)$$

and for  $x \in N$ ,

$$\dot{m}_g(x) = \dot{m}_g(0) + \dot{m}_l(0) - \dot{m}_l(x)$$

and for  $v_g$ , the initial condition and general expression for  $x \in N$  is

$$v_g(0) = \frac{\dot{m}_g(0)R_UT_g(0)}{MW_g(\phi_g(0))PA_{cc}} \quad v_g(x) = \frac{\dot{m}_g(x)R_UT_g(x)}{MW_g(\phi_g(x))PA_{cc}}$$

Taking stock of everything so far, the overall outlook is the following:

- chamber parameters are fixed and given

$$(A_{\text{tot fuel inj}}, A_{cc}, l_{cc}) \in (\mathbb{R}^+)^3$$

- inlet parameters

$$((T_{\text{in}}, P, \phi_g(0)), \phi_{\text{overall}}, v_d(0), D(0)) \in (\Sigma_{x=0}, C^\infty(\Sigma_{x=0}), TN, C^\infty(N))$$

The use of the Clausius-Clapeyron relation for the vapor pressure, and how the thermodynamic system straddles between the liquid phase and vapor (gas) phase, yields a one-to-one mapping (isomorphism): from  $p = p_0 \exp \left( L \left( \frac{1}{T} - \frac{1}{T^o} \right) \right)$ ,

$$p \xrightarrow{\text{Cl} - \text{Cl}} T(p) = T_{\text{boil}}$$

So for

$$\dot{m}_l(0) = \rho_l(T_b, P)v_d(0)A_{\text{tot fuel inj}}$$

$$\dot{m}_{\text{Ox}}(0) = \frac{\dot{m}_l(0)}{(\phi_{\text{overal}} - \phi_g(0))\left(\frac{F}{O}\right)_{\Phi=1}}$$

$$\dot{m}_g(0) = \left( \phi_{\text{overall}} \left( \frac{F}{O} \right)_{\Phi=1} + 1 \right) \dot{m}_{\text{Ox}}(0) - \dot{m}_l(0)$$

$$v_g(0) = \frac{\dot{m}_g(0)R_UT_g(0)}{MW_g(\phi_g(0))PA_{cc}}$$

and then for general  $x \in N$ ,

$$\dot{m}_l(x) = \dot{m}_l(0) \left( \frac{D}{D_0} \right)^3$$

$$\dot{m}_g(x) = \dot{m}_g(0) + \dot{m}_l(0) - \dot{m}_l(x)$$

$$\phi_g(x) = \frac{\frac{\dot{m}_g(x)}{\dot{m}_{\text{Ox}}(0)} - 1}{\left(\frac{F}{O}\right)_{\Phi=1}}$$

$$v_g(x) = \frac{\dot{m}_g(x)R_UT_g(x)}{MW_g(\phi_g(x))PA_{cc}}$$

$$\text{chamber parameters, inlet conditions} \longmapsto (\dot{m}_l(0), \dot{m}_{\text{Ox}}(0), \dot{m}_g(0), v_g(0))$$

$$\text{chamber parameters, inlet conditions} \times (\dot{m}_l(0), \dot{m}_{\text{Ox}}(0), \dot{m}_g(0)) \xrightarrow{(T_g(x), D(x), v_d(x))} (\dot{m}_l(x), \dot{m}_g(x)) \times (\phi_g(x), v_g(x))$$

### Part 7. Numerical Computation; Scientific Computation

#### 24. INTERPOLATION AND EXTRAPOLATION

cf. Ch. 3, Interpolation and Extrapolation, Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#).

If we know  $f(x)$  at  $x_0, x_1, \dots, x_{N-1}$ ,  $x_0 < x_1 < \dots < x_{N-1}$ ,

2 interpolation processes:

- (1) Fit (once) interpolating function to data pts. provided
- (2) Evaluate (as many tiumes as you wish) that interpolating function at target pt.  $x$

Cons: typically less computationally efficient, more susceptible to round off error.

- (1) Find right starting position in table ( $x_i$  or  $i$ )
- (2) Perform interpolation (construct functional estimate  $f(x)$ ) using  $M$  nearby values (e.g. centered on  $M \ll N$  and  $O(M^2)$  time (operations))

local interpolation using  $M$  nearest neighbor points gives interpolated values  $f(x)$  that don't in general, have cont. 1st. or higher derivatives,

because, as  $x = x_i$ ,  $x = x_j$  tabulated values, interpolation scheme switches which tabulated pts. are "local"

where continuity of derivatives is concern, use "stiffer" interpolation by *spline* function. *spline* is a polynomial between each pair of table pts., but coefficients determined "slightly" nonlocally,

- nonlocality to guarantee global smoothness in interpolated function up to some order of derivative.

*order* of the interpolation =  $M - 1$ ,  $M \equiv$  number of pts. used in interpolation scheme.

abscissas  $x_j$ ,  $j = 0, \dots, N - 1$

abscissas either monotonically increasing or monotonically decreasing.

Given  $M \leq N$ , number  $x$ ,

find  $j_{\text{lo}} \in \mathbb{Z}$  s.t.  $x$  centered among  $M$  abscissas  $x_{j_{\text{lo}}}, \dots, x_{j_{\text{lo}}}, \dots, x_{j_{\text{lo}}+M-1}$ , i.e.  $x \in [x_m, x_{m+1}]$ , where

$$(109) \quad m = j_{\text{lo}} + \lfloor \frac{M-2}{2} \rfloor$$

and  $j_{\text{lo}} \not\prec 0$  and  $j_{\text{lo}} + M - 1 \not\succ N - 1$

If  $M \bmod 2 = 0$ ,  $M = 2N$ , so  $\frac{M-2}{2} = N - 1$ .

If  $M \bmod 2 = 1$ ,  $M = 2N + 1$ , so  $\frac{M-2}{2} = \frac{2(N-1)+1}{2} \mapsto N - 1$ , if  $M > 2$ .

#### 24.1. Polynomial Interpolation and Extrapolation.

24.1.1. *Lagrange polynomials or Lagrange's classical formula.* Given  $M$  pts.

$$(x_0, y_0), (x_1, y_1), \dots, (x_{M-1}, y_{M-1})$$

where no 2  $x_j$ 's,  $j = 0, 1 \dots M - 1$  are equal, interpolation polynomial in Lagrange form is a linear combination:

$$(110) \quad L(x) := \sum_{j=0}^{M-1} y_j l_j(x)$$

of Lagrange basis polynomials  $l_j(x)$ :

$$(111) \quad l_j(x) := \prod_{\substack{0 \leq m \leq M-1 \\ m \neq j}} \frac{x - x_m}{x_j - x_m} = \frac{(x - x_0)}{(x_j - x_0)} \cdots \frac{(x - x_{j-1})(x - x_{j+1})}{(x_j - x_{j-1})(x_j - x_{j+1})} \cdots \frac{(x - x_{M-1})}{(x_j - x_{M-1})}$$

and so, for change of notation,

$$L(x) \equiv P(x) = \frac{(x - x_1)(x - x_2) \dots (x - x_{M-1})}{(x_0 - x_1)(x_0 - x_2) \dots (x_0 - x_{M-1})} y_0 + \frac{(x - x_0)(x - x_1) \dots (x - x_{M-2})}{(x_1 - x_0)(x_1 - x_2) \dots (x_1 - x_{M-1})} y_1 + \dots$$

$$+ \frac{(x - x_0)(x - x_1) \dots (x - x_{M-2})}{(x_{M-1} - x_0)(x_{M-1} - x_1) \dots (x_{M-1} - x_{M-2})} y_{M-1}$$

For  $M = 2, M = 3$  cases,

$$L^{(2)}(x) = \frac{x - x_1}{x_0 - x_1} y_0 + \frac{x - x_0}{x_1 - x_0} y_1$$

$$L^{(3)}(x) = \frac{(x - x_1)(x - x_2)}{(x_0 - x_1)(x_0 - x_2)} y_0 + \frac{(x - x_0)(x - x_2)}{(x_1 - x_0)(x_1 - x_2)} y_1 + \frac{(x - x_0)(x - x_1)}{(x_2 - x_0)(x_2 - x_1)} y_2$$

### 25. INTEGRATION OF ODES: RUNGE-KUTTA METHODS

#### 25.1. Runge Kutta Methods. cf. [Appendix A, Runge-Kutta Methods](#)

Consider an initial value problem (IVP):

$$(112) \quad \frac{d\mathbf{x}}{dt} = \mathbf{f}(t, \mathbf{x}(t))$$

where

$$\mathbf{x}(t) = (x_1(t), x_2(t), \dots, x_n(t))$$

$$f \in [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n$$

$$\mathbf{x}(0) = \mathbf{x}_0$$

Divide  $[a, b]$  into  $M$  equal subintervals, select mesh points  $t_j$ :

$$t_j = a + jh, \quad j = 0, 1 \dots M, \quad h = \frac{b - a}{M}$$

Family of explicit Runge-Kutta (RK) methods of  $m$ th stage given by

$$(113) \quad \begin{aligned} \mathbf{x}(t_{n+1}) &\equiv \mathbf{x}_{n+1} = \mathbf{x}_n + h \sum_{i=1}^m c_i k_i \text{ where} \\ k_1 &= \mathbf{f}(t_n, \mathbf{x}_n) \\ k_2 &= \mathbf{f}(t_n + \alpha_2 h, \mathbf{x}_n + h\beta_{21}k_1(t_n, \mathbf{x}_n)) \\ k_3 &= \mathbf{f}(t_n + \alpha_3 h, \mathbf{x}_n + h(\beta_{31}k_1(t_n, \mathbf{x}_n) + \beta_{32}k_2(t_n, \mathbf{x}_n))) \\ &\vdots \\ k_m &= \mathbf{f}(t_n + \alpha_m h, \mathbf{x}_n + h \sum_{j=1}^{m-1} \beta_{mj} k_j) \end{aligned}$$

These data are usually arranged in a *Butcher tableau*:

0					
$\alpha_2$	$\beta_{21}$				
$\alpha_3$	$\beta_{31}$	$\beta_{32}$			
$\vdots$				$\ddots$	
$\alpha_m$	$\beta_{m1}$	$\beta_{m2}$	$\dots$	$\beta_{m,m-1}$	
	$c_1$	$c_2$	$\dots$	$c_{m-1}$	$c_m$

In the notation of Hairer, Nørsett, and Wanner (1993) [37], Eq. (1.8'),

$$(114) \quad \begin{array}{c|ccccc} 0 & & & & & \\ c_2 & a_{21} & & & & \\ c_3 & a_{31} & a_{32} & & & \\ \vdots & & & & \ddots & \\ c_s & a_{s1} & a_{s2} & \dots & a_{s,s-1} & \\ \hline & b_1 & b_2 & \dots & b_{s-1} & b_s \end{array}$$

To specify a method, we need,

$$(115) \quad \begin{aligned} m &\equiv \text{number of stages,} \\ \alpha_i &\quad (i = 2, 3 \dots m) \\ \beta_{ij} &\quad (1 \leq j < i \leq m) \text{ and} \\ c_i &\quad (i = 1, 2, \dots m) \end{aligned}$$

In summary,  
to solve the following initial value problem (IVP):

$$(116) \quad \boxed{\frac{d\mathbf{x}}{dt} = \mathbf{f}(t, \mathbf{x}(t)), \quad \mathbf{x} \in \mathbb{R}^n, \quad \mathbf{f} \in [a, b] \times \mathbb{R}^n \rightarrow \mathbb{R}^n}$$

We need

$$(117) \quad \begin{aligned} \alpha_i &\quad (i = 2, 3, \dots m) \\ \beta_{ij} &\quad (1 \leq j < i \leq m), \quad (m(m-1) \text{ total } \beta_{ij} \text{ coefficients}) \\ c_i &\quad (i = 1, 2, \dots m) \end{aligned}$$

To calculate, for

$$(118) \quad \boxed{t_j = a + jh, \quad j = 0, 1, \dots M, \quad h = \frac{b - a}{M}}$$

So for  $n \in 0, 1, \dots M$

(119)

$$k_1 = \mathbf{f}(t_n, \mathbf{x}_n)$$
$$k_l = \mathbf{f}(t_n + \alpha_l h, \mathbf{x}_n + h \sum_{j=1}^{l-1} \beta_{lj} k_j), \quad l = 2, \dots m$$

To obtain the result

(120)

$$\mathbf{x}(t_{n+1}) \equiv \mathbf{x}_{n+1} = \mathbf{x}_n + h \sum_{i=1}^m c_i k_i$$

Compare this notation, include Eq. [113](#) to what’s in Hairer, Nørsett, and Wanner (1993) [\[37\]](#):

**Definition 7** (1.1, Hairer, Nørsett, and Wanner (1993) [\[37\]](#)). *Let  $s \in \mathbb{Z}$ ,  $s \equiv$  number of stages, and*

(121)

$$a_{21}, a_{31}, a_{32}, \dots a_{s1}, a_{s2}, \dots a_{s,s-1},$$
$$b_1, \dots b_s,$$
$$c_2, \dots c_s$$

be real coefficients.

Then

(122)

$$k_1 = f(x_0, y_0)$$
$$k_2 = f(x_0 + c_2 h, y_0 + h a_{21} k_1)$$
$$k_3 = f(x_0 + c_3 h, y_0 + h(a_{31} k_1 + a_{32} k_2))$$
$$\dots$$
$$k_s = f(x_0 + c_s h, y_0 + h(a_{s1} k_1 + \dots + a_{s,s-1} k_{s-1}))$$
$$y_1 = y_0 + h(b_1 k_1 + \dots + b_s k_s)$$

is called a  $s$ -stage explicit Runge-Kutta method (ERK) for Eq. [\(1.1\)](#), Hairer, Nørsett, and Wanner (1993) [\[37\]](#), i.e.  $y' = f(x, y)$ ,  $y(x_0) = y_0$ .

Usually,  $c_i$  satisfy conditions:

$$c_2 = a_{21}$$
$$c_3 = a_{31} + a_{32}$$
$$\dots$$
$$c_s = a_{s1} + \dots + a_{s,s-1}$$

or briefly

$$c_i = \sum_{j=1}^{i-1} a_{ij}$$

The notation I will use is the following, replacing Eq. [122](#):

$$\mathbf{y}(x_n) \equiv \mathbf{y}_n = \mathbf{y}_{n-1} + h(b_1 \mathbf{k}_1 + \dots + b_s \mathbf{k}_s) = \mathbf{y}_{n-1} + h \sum_{i=1}^s b_i \mathbf{k}_i$$

(123)

$$\mathbf{k}_1 = \mathbf{f}(x_n, \mathbf{y}_n)$$
$$\mathbf{k}_2 = \mathbf{f}(x_n + c_2 h, \mathbf{y}_n + h a_{21} \mathbf{k}_1)$$
$$\mathbf{k}_3 = \mathbf{f}(x_n + c_3 h, \mathbf{y}_n + h(a_{31} \mathbf{k}_1 + a_{32} \mathbf{k}_2))$$
$$\dots$$
$$\mathbf{k}_s = \mathbf{f}(x_n + c_s h, \mathbf{y}_n + h(a_{s1} \mathbf{k}_1 + \dots + a_{s,s-1} \mathbf{k}_{s-1})) = \mathbf{f}(x_n + c_s h, \mathbf{y}_n + h \sum_{j=1}^{s-1} a_{sj} \mathbf{k}_j)$$

**Definition 8** (1.2, Hairer, Nørsett, and Wanner (1993) [\[37\]](#)). *Runge-Kutta method (1.8, Hairer, Nørsett, and Wanner (1993) [\[37\]](#), [122](#)) has order  $p$  if for sufficiently smooth problems (1.1, Hairer, Nørsett, and Wanner (1993) [\[37\]](#))*

(124)

$$\|y(x_0 + h) - y_1\| \leq K h^{p+1}$$

cf. Eqn. (1.10) of Hairer, Nørsett, and Wanner (1993) [\[37\]](#).  
i.e. if Taylor series for exact solution  $y(x_0 + h)$  and for  $y_1$  coincide up to (and including) term  $h^p$ .  
For notation, rewrite Eq. [124](#) as

(125)

$$\|\mathbf{y}(x_n + h) - \mathbf{y}_n\| \leq K h^{p+1}$$

where

(126)

$$\mathbf{y}_n \equiv \mathbf{y}(x_n) = \mathbf{y}_{n-1} + h \sum_{i=1}^s b_i \mathbf{k}_i$$

25.1.1. *Examples.* 1.  $m = 1$ . Then

$$k_1 = \mathbf{f}(t_n, \mathbf{x}_n)$$
$$\mathbf{x}_{n+1} = \mathbf{x}_n + h c_1 k_1$$

By Taylor expansion,

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h \dot{\mathbf{x}}|_{t_n} + \dots = \mathbf{x}_n + h \mathbf{f}(t_n, \mathbf{x}_n) + \mathcal{O}(h^2)$$

So  $\mathbf{x}_{n+1} = \mathbf{x}_n + h f(t_n, \mathbf{x}_n)$ .

$$c_1 = 1$$

Using the notation of Hairer, Nørsett, and Wanner (1993) [\[37\]](#), the case of  $s = 1$  is the Euler case:  
 $s = 1 \implies b_1$

$$\mathbf{y}(x_n) = \mathbf{y}(x_{n-1}) + h b_1 \mathbf{k}_1 = \mathbf{y}(x_{n-1}) + h b_1 \mathbf{f}(x_n, \mathbf{y}_n)$$
$$x_n = a + n h, h = \frac{b-a}{M}, \quad n = 0, 1, \dots M$$

$$m = 2$$



Now consider the Taylor series expansion

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h \dot{\mathbf{x}}|_{t_n} + \frac{h^2}{2} \frac{d^2 \mathbf{x}}{dt^2} \Big|_{t_n} + \frac{h^3}{6} \frac{d^3 \mathbf{x}}{dt^3} \Big|_{t_n} + \frac{h^4}{24} \frac{d^4 \mathbf{x}}{dt^4} \Big|_{t_n} + \mathcal{O}(h^5)$$

So

$$\begin{aligned} \frac{d^3 \mathbf{x}}{dt^3} &= \frac{\partial^2 f}{\partial t^2} + f_i \frac{\partial^2 f}{\partial x_i \partial t} + \frac{\partial f_i}{\partial t} \frac{\partial f}{\partial x_i} + \frac{\partial f_i}{\partial x_j} f_j \frac{\partial f}{\partial x_i} \\ \frac{d^4 \mathbf{x}}{dt^4} &= f_{ttt} + f_{x_i tt} f_i + f_{t,i} f_{x_i t} + f_{x_j, i} f_j f_{x_i t} + f_i f_{x_i tt} + f_i f_j f_{x_i t x_j} + f_{tt, i} f_{x_i} + f_{t x_j, i} f_j f_{x_i} + f_{t, i} f_{x_i t} + f_{t, i} f_j f_{x_i t x_j} + \\ &\quad + \left( \frac{\partial}{\partial t} + \nabla \right) \left( \frac{\partial f_i}{\partial x_j} f_j \frac{\partial f}{\partial x_i} \right) \end{aligned}$$

TODO: Complete proof.

So

$$\begin{aligned} c_1 + c_2 + c_3 + c_4 &= 1 \\ \beta_{21} &= \alpha_2 \\ \beta_{31} + \beta_{32} &= \alpha_3 \\ c_2 \alpha_2 + c_3 \alpha_3 + c_4 \alpha_4 &= \frac{1}{2} \\ c_2 \alpha_2^2 + c_3 \alpha_3^2 + c_4 \alpha_4^2 &= \frac{1}{3} \\ c_2 \alpha_2^3 + c_3 \alpha_3^3 + c_4 \alpha_4^3 &= \frac{1}{4} \\ c_3 \alpha_2 \beta_{32} + c_4 (\alpha_2 \beta_{42} + \alpha_3 \beta_{43}) &= \frac{1}{6} \\ c_3 \alpha_2 \alpha_3 \beta_{32} + c_4 \alpha_4 (\alpha_2 \beta_{42} + \alpha_3 \beta_{43}) &= \frac{1}{8} \\ c_3 \alpha_2^2 \beta_{32} + c_4 (\alpha_2^2 \beta_{42} + \alpha_3^2 \beta_{43}) &= \frac{1}{12} \\ c_4 \alpha_2 \beta_{32} \beta_{43} &= \frac{1}{24} \end{aligned} \tag{128}$$

13 unknowns with 11 equations yields 2 additional conditions.

For

$$\alpha_2 = \frac{1}{2}, \quad \beta_{31} = 0$$

We obtain the classical RK4 method:

In Butcher tableau form:

$$\begin{array}{c|ccc} 0 & & & \\ \frac{1}{2} & \frac{1}{2} & & \\ \frac{1}{2} & 0 & \frac{1}{2} & \\ 1 & 0 & 0 & 1 \\ & \frac{1}{6} & \frac{1}{3} & \frac{1}{3} & \frac{1}{6} \end{array}$$

$$k_1 = \mathbf{f}(t_n, \mathbf{x}_n)$$

$$k_2 = \mathbf{f}(t_n + \alpha_2 h, \mathbf{x}_n + h \beta_{21} k_1)$$

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h c_1 k_1 + h c_2 k_2$$

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h \dot{\mathbf{x}}|_{t_n} + \frac{h^2}{2} \frac{d^2 \mathbf{x}}{dt^2} \Big|_{t_n} + \mathcal{O}(h^3)$$

$$\begin{aligned} \frac{d^2 \mathbf{x}}{dt^2} &= \frac{d}{dt} \mathbf{f}(t, \mathbf{x}(t)) = \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{x}) + \dot{\mathbf{x}}(t) \frac{\partial \mathbf{f}}{\partial \mathbf{x}} = \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{x}) + f_i(t, \mathbf{x}) \frac{\partial f}{\partial x_i} \\ \implies \mathbf{x}_{n+1} - \mathbf{x}_n &= h \mathbf{f}(t_n, \mathbf{x}_n) + \frac{h^2}{2} \left( \frac{\partial}{\partial t} \mathbf{f}(t, \mathbf{x}) + f_i(t_n, \mathbf{x}_n) \frac{\partial f}{\partial x_i}(t_n, \mathbf{x}_n) \right) \end{aligned}$$

$$k_2 = f(t_n + \alpha_2 h, \mathbf{x}_n + h \beta_{21} k_1) = f(t_n, \mathbf{x}_n) + \alpha_2 h \frac{\partial}{\partial t} \mathbf{f}(t_n, \mathbf{x}_n) + h \beta_{21} k_{1,i} \frac{\partial}{\partial x_i} \mathbf{f}(t_n, \mathbf{x}_n) + \mathcal{O}(h^2)$$

$$h f(t_n, \mathbf{x}_n) = (h c_1 + h c_2) k_1 \implies c_1 + c_2 = 1$$

$$\frac{h^2}{2} = h^2 c_2 \alpha_2 \implies c_2 \alpha_2 = \frac{1}{2}$$

$$\frac{h^2}{2} = h^2 c_2 \beta_{21} \implies c_2 \beta_{21} = \frac{1}{2}$$

Let  $\alpha_2 = 1$ . Then  $c_2 = \frac{1}{2}$ ,  $c_1 = \frac{1}{2}$ .

$$\implies \mathbf{x}_{n+1} = \mathbf{x}_n + \frac{h}{2} (\mathbf{f}(t_n, \mathbf{x}_n) + \mathbf{f}(t_n + h, \mathbf{x}_n + h \mathbf{f}(t_n, \mathbf{x}_n)))$$

This is called Heun's method.

Let  $\alpha_2 = \frac{1}{2}$ . Then  $c_2 = 1$ ,  $\beta_{21} = \frac{1}{2}$ ,  $c_1 = 0$

$$\mathbf{x}_{n+1} = \mathbf{x}_n + h \mathbf{f}(t_n + \frac{h}{2}, \mathbf{x}_n + \frac{h}{2} \mathbf{f}(t_n, \mathbf{x}_n))$$

This is called the RK2 method.

Using the notation of Hairer, Nørsett, and Wanner (1993) [37],

$s = 2 \implies$

$a_{21}$

$b_1, b_2$

$c_2$

$$\mathbf{k}_1 = \mathbf{f}(x_n, \mathbf{y}_n)$$

$$\mathbf{k}_2 = \mathbf{f}(x_n + c_2 h, \mathbf{y}_n + h a_{21} \mathbf{k}_1)$$

and

$$\mathbf{y}(x_n) = \mathbf{y}_{n-1} + h(b_1 \mathbf{k}_1 + b_2 \mathbf{k}_2)$$

25.1.2. *Runge-Kutta 4, RK4 Methods.*  $m = 4$ .

By matching coefficients with Taylor series, then for

$$\begin{aligned} k_1 &= \mathbf{f}(t_n, \mathbf{x}_n) \\ k_2 &= \mathbf{f}(t_n + \alpha_2 h, \mathbf{x}_n + h \beta_{21} k_1(t_n, \mathbf{x}_n)) \\ k_3 &= \mathbf{f}(t_n + \alpha_3 h, \mathbf{x}_n + h(\beta_{31} k_1(t_n, \mathbf{x}_n) + \beta_{32} k_2(t_n, \mathbf{x}_n))) \\ k_4 &= \mathbf{f}(t_n + \alpha_4 h, \mathbf{x}_n + h(\beta_{41} k_1 + \beta_{42} k_2 + \beta_{43} k_3)) \end{aligned} \tag{127}$$

Or

$$\begin{aligned}\alpha_2 &= \alpha_3 = \frac{1}{2}, & \alpha_4 &= 1 \\ \beta_{21} &= \frac{1}{2} \\ \beta_{31} &= 0, & \beta_{32} &= \frac{1}{2} \\ \beta_{41} &= \beta_{42} = 0, & \beta_{43} &= 1 \\ c_1 &= \frac{1}{6}, & c_2 &= c_3 = \frac{1}{3}, & c_4 &= \frac{1}{6}\end{aligned}$$

If you use our notation, in Eq. [123](#),

$$\begin{aligned}k_1 &= f(x_n, y_n) \\ k_2 &= f(x_n + c_2 h, y_n + h a_{21} k_1) \\ k_3 &= f(x_n + c_3 h, y_n + h(a_{31} k_1 + a_{32} k_2)) \\ k_4 &= f(x_n + c_4 h, y_n + h(a_{41} k_1 + a_{42} k_2 + a_{43} k_3)) \\ a_{21} &= \frac{1}{2} \\ a_{32} &= \frac{1}{2} & a_{31} &= 0 \\ a_{43} &= 1 & a_{42} &= a_{41} = 0 \\ c_2 &= \frac{1}{2} & c_3 &= \frac{1}{2} & c_4 &= 1 \\ b_1 &= \frac{1}{6} & b_2 &= b_3 = \frac{1}{3} & b_4 &= \frac{1}{6} \\ k_1 &= f(x_n, y_n) \\ k_2 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2} k_1) \\ k_3 &= f(x_n + \frac{h}{2}, y_n + \frac{h}{2} k_2) \\ k_4 &= f(x_n + h, y_n + h k_3) \\ y_n &= y_{n-1} + h \sum_{i=1}^4 b_i k_i\end{aligned}\tag{129}$$

25.1.3. *Explicit Runge-Kutta Methods.* Recall Eq. [123](#):

$$\begin{aligned}\mathbf{k}_i &= \mathbf{f}(x_n + c_i h, \mathbf{y}_n + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j) \\ \mathbf{y}(x_{n+1}) &= \mathbf{y}_n + h \sum_{i=1}^s b_i \mathbf{k}_i\end{aligned}\tag{130}$$

This can be rewritten as

$$\begin{aligned}g_i &= y_m + hJ \sum_{j=1}^{i-1} a_{ij} g_j \\ y_{m+1} &= y_m + hJ \sum_{j=1}^s b_j g_j\end{aligned}\tag{131}$$

or

$$\begin{aligned}\mathbf{g}_i &= \mathbf{y}_n + hJ \sum_{j=1}^{i-1} a_{ij} \mathbf{g}_j \\ \mathbf{y}_{n+1} &= \mathbf{y}_n + hJ \sum_{j=1}^s b_j \mathbf{g}_j\end{aligned}\tag{132}$$

cf. Eq. (2.7) of Hairer and Wanner (2010) [\[38\]](#).

Insert  $\mathbf{g}_j$  into the expression for  $\mathbf{y}_{n+1}$ , so then

$$\mathbf{y}_{n+1} = R(hJ) \mathbf{y}_m\tag{133}$$

where

$$R(z) = 1 + z \sum_j b_j + z^2 \sum_{j,k} b_j a_{jk} + z^3 \sum_{j,k,l} b_j a_{jk} a_{kl} + \dots\tag{134}$$

is a polynomial of degree  $\leq s$ . This is Eq. (2.8) of Hairer and Wanner (2010) [\[38\]](#).

**25.2. Order Conditions for Runge-Kutta Methods.** cf. pp. 143, II.2 Order Conditions for Runge-Kutta Methods, Hairer, Nørsett, and Wanner (1993) [\[37\]](#)

Replace  $\mathbf{k}_i$  by  $\mathbf{g}_i$  s.t.  $\mathbf{k}_i = \mathbf{f}(\mathbf{g}_i)$

**25.3. Adaptive Stepsize Control, Stepsize Selection for Runge-Kutta and Practical Error Estimation.**

25.3.1. *Richardson Extrapolation.* Suppose that, given initial value  $(x_0, y_0)$ , step size  $h$ , compute 2 steps (step doubling).

Compute, starting from  $(x_0, y_0)$  with one big step  $2h$ , approximate solution  $y_1$ , and approximate solution  $y_2$  (2 steps of size  $h$ ).

Let exact solution be  $y(x_0 + 2h) \equiv y(x + 2h)$ :

$$e_1 = y(x_0 + h) - y_1 = Ch^{p+1} + \mathcal{O}(h^{p+2})\tag{135}$$

cf. Eq. 4.1, Hairer, Nørsett, and Wanner (1993) [\[37\]](#), or in my notation,

$$e_1 \equiv \text{err}_1 = y(x + 2h) - y_1 = C(2h)^{p+1} + \mathcal{O}(h^{p+2})\tag{136}$$

with  $C = 1$ ,  $p = 4$  as an example.

$C$  contains error coefficients of the method and elementary differentials  $F^J(t)(y_0) \equiv f^J(x)(y_0)$  of order  $p + 1 = 5$ . Error of  $y_2$  consists of 2 parts:

- transported error of first step,  $(I + h \frac{\partial f}{\partial y} + \mathcal{O}(h^2))e_1$
- local error of 2nd. step

$$e_1 = y(x_0 + h) - y_1 = Ch^{p+1} + \mathcal{O}(h^{p+2}) \text{ s.t. } y_1 = y_0 + \mathcal{O}(h)$$

Thus,

$$\begin{aligned} \implies e_2 &= y(x_0 + 2h) - y_2 = (I + \mathcal{O}(h))Ch^{p+1} + (C + \mathcal{O}(h))h^{p+1} + \mathcal{O}(h^{p+2}) = \\ &= 2Ch^{p+1} + \mathcal{O}(h^{p+2}) \\ \implies y(x_0 + 2h) - y &= C(2h)^{p+1} + \mathcal{O}(h^{p+2}) \end{aligned}$$

25.3.2. *Embedded Runge-Kutta Formulas.* For  $s = 5$ , and comparing with Eq. [123](#), fifth order Runge-Kutta formula is, in Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#)’s notation,

$$\begin{aligned} k_1 &= hf(x_n, y_n) \\ k_2 &= hf(x_n + c_2h, y_n + a_{21}k_1) \\ &\dots \\ k_6 &= hf(x_n + c_6h, y_n + \sum_{j=1}^5 a_{6j}k_j) \\ y_{n+1} &= y_n + \sum_{i=1}^6 b_i k_i + \mathcal{O}(h^6) \end{aligned} \tag{137}$$

cf. Eq. (17.2.4), pp. 912 of Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#).

The *embedded* 4th-order formula is of general form

$$y_{n+1}^* = y_n + \sum_{i=1}^6 b_i^* k_i + \mathcal{O}(h^5) \tag{138}$$

cf. pp. 912, Eq. (17.2.5), Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#), so that (the general form of) the error estimate is

$$\Delta \equiv y_{n+1} - y_{n+1}^* = \sum_{i=1}^6 (b_i - b_i^*) k_i \tag{139}$$

cf. pp. 921, Eq. (17.2.6), Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#).

Note the notation change from Eq. [138](#) to Eqns. 4.7, 4.7’ in Hairer, Nørsett, and Wanner (1993) [\[37\]](#):

$$\begin{aligned} y_{n+1} &\equiv y_1 = y_0 + h \sum_{i=1}^s b_i k_i \\ y_{n+1}^* &\equiv \hat{y}_1 = y_0 + h \sum_{s=1}^s \hat{b}_i k_i \end{aligned}$$

The Butcher tableau also changes: in the notation of Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#) on pp.913 for Dormand-Prince 5(4) Parameters for Embedded Runge-Kutta Method, it’s implied that the table becomes

0					
$c_2$	$a_{21}$				
$c_3$	$a_{31}$	$a_{32}$			
$\vdots$					$\ddots$
$c_s$	$a_{s1}$	$a_{s2}$	$\dots$	$a_{s,s-1}$	
	$b_1$	$b_2$	$\dots$	$b_{s-1}$	$b_s$
	$b_1^*$	$b_2^*$	$\dots$	$b_{s-1}^*$	$b_s^*$

In the notation of Hairer, Nørsett, and Wanner (1993) [\[37\]](#), Eq. (4.6), pp. 166, II Runge-Kutta and Extrapolation Methods,

$$\begin{array}{c|cccccc} 0 & & & & & & \\ c_2 & a_{21} & & & & & \\ c_3 & a_{31} & a_{32} & & & & \\ \vdots & & & & \ddots & & \\ c_s & a_{s1} & a_{s2} & \dots & a_{s,s-1} & & \\ \hline & b_1 & b_2 & \dots & b_{s-1} & b_s & \\ \hline & \widehat{b}_1 & \widehat{b}_2 & \dots & \widehat{b}_{s-1} & \widehat{b}_s & \end{array} \tag{140}$$

Consider that the use of  $y_{n+1}$  itself to provide a 7th ( $s + 1$ th) stage. Because  $f(x_n + h, y_{n+1})$  has to be evaluated anyway to start the next step, this costs nothing. This trick is called FSAL (First Same As Last).

cf. pp. 912-913, 17.2 Adaptive Stepsize Control for Runge-Kutta, Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#) cf. pp. 166-167, II.4 Practical Error Estimation and Stepsize Selection, Hairer, Nørsett, and Wanner (1993) [\[37\]](#)

So then

$$k_7 = hf(x_n + h, y_n + \sum_{i=1}^6 b_i k_i) = hf(x_n + h, y_{n+1})$$

$$c_7 = 1$$

$$k_{s+1} = f(x_n + h, y_n + h \sum_{i=1}^s b_i k_i) = f(x_n + c_{s+1}h, y_n + h \sum_{i=1}^s a_{s+1,i} k_i) = f(x_n + h, y_{n+1})$$

25.3.3. *Automatic Stepsize Control.* D’ordinaire, on se content de multiplier ou de diviser par 2 la valeur du pas. -Ceschino 1961, pp. 167 II. Runge-Kutta and Extrapolation Methods of Hairer, Nørsett, and Wanner (1993) [\[37\]](#) (Usually, we’re content multiplying or dividing by the value 2 each step).

Comparing Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#) and Hairer, Nørsett, and Wanner (1993) [\[37\]](#), respectively,

$$|\Delta| = |y_{n+1} - y_{n+1}^*| \leq \mathbf{scale} \tag{141}$$

cf. Eq. (17.2.7), Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#), or

$$|\Delta| = \|y_{n+1} - y_{n+1}^*\| \leq \mathbf{sc}_i \tag{142}$$

cf. Eq. (4.10), Hairer, Nørsett, and Wanner (1993) [\[37\]](#)

Furthermore,

$$\mathbf{scale} = \mathbf{atol} + |y| \mathbf{rtol} \tag{143}$$

cf. Eq. (17.28), Press, Teukolsky, Vetterling, Flannery (2007) [\[36\]](#), with  $|y| \mapsto \max(|y_n|, |y_{n+1}|)$ .

In Hairer, Nørsett, and Wanner (1993) [\[37\]](#)’s notation,

$$\mathbf{sc}_i = \mathbf{Atol}_i + \max(|y_{0i}|, |y_{1i}|) \cdot \mathbf{Rtol}_i \tag{144}$$

cf. Eq. (4.10), Hairer, Nørsett, and Wanner (1993) [\[37\]](#), where  $\mathbf{Atol}_i \equiv \mathbf{atol}$ ,  $\mathbf{Rtol}_i \equiv \mathbf{rtol}$  are absolute error tolerance, relative error tolerance, respectively.

Usually  $\mathbf{Atol}_i$ ,  $\mathbf{rtol}_i$  are both different from 0; they may also depend on the component of the solution, says Hairer, Nørsett, and Wanner (1993) [\[37\]](#) (pp. 167), but Press, Teukolsky, Vetterling, Flannery (2007) says that while  $\mathbf{atol}, \mathbf{rtol}$  could be different for each component of  $\mathbf{y}$ , we’ll take them as constant.

Take as a measure of error Eq. (4.11) of Hairer, Nørsett, and Wanner (1993) [\[37\]](#)

$$\text{err} = \sqrt{\frac{1}{n} \sum_{i=1}^n \left( \frac{y_{1i} - \hat{y}_{1i}}{\mathbf{sc}_i} \right)^2} \tag{145}$$

and compare this to Press, Teukolsky, Vetterling, Flannery (2007) for Eq. (17.2.9)

$$(146) \quad \mathbf{err} = \sqrt{\frac{1}{N} \sum_{i=0}^{N-1} \left( \frac{\Delta_i}{\mathbf{scale}_i} \right)^2}$$

Consider Eq. 145,  $\mathbf{err} = \sqrt{\frac{1}{N} \sum_{i=0}^{N-1} \left( \frac{\Delta_i}{\mathbf{sc}_i} \right)^2}$  where  $\Delta_i \equiv y_{n+1,i} - y_{n+1,i}^*$ .

Since  $\mathbf{y}_1 \equiv \mathbf{y}_{n+1}$  is of order  $p$ ,  $\hat{y}_1 \equiv \mathbf{y}_{n+1}$  of order  $\hat{p}$  (usually  $\hat{p} = p - 1$ , or  $\hat{p} = p + 1$ ), with  $\hat{p} = 5$ .

From Eqns. (17.2.4), (17.2.5) of Press, Teukolsky, Vetterling, Flannery (2007) [36], i.e. Eq. 123, for the embedded 4th order Runge-Kutta method,  $\Delta_i$  scales as  $h^5 \equiv h^{\hat{p}} = h^{p+1}$ , and hence so does  $\mathbf{err}$ .

If we take step  $h_1$  and produce error  $\mathbf{err}_1$ , then for step  $h_0$ , that would've given some other value  $\mathbf{err}_0$ , as explained on pp. 913 of Press, Teukolsky, Vetterling, Flannery (2007) [36], the error  $\mathbf{err}_0$  is estimated as

$$(147) \quad h_0 = h_1 \left| \frac{\mathbf{err}_0}{\mathbf{err}_1} \right|^{1/5} \equiv h_1 \left| \frac{\mathbf{err}_0}{\mathbf{err}_1} \right|^{\frac{1}{q+1}}$$

cf. (17.2.10) of Press, Teukolsky, Vetterling, Flannery (2007) [36], pp. 913.

Or from error behavior  $\mathbf{err} \approx C \cdot h^{q+1}$ , and from  $1 \approx C \cdot h_{\text{opt}}^{q+1}$  (where  $q = \min(p, \hat{p})$ ),  $\mathbf{err}_0 = 1$  is "an efficient integration", i.e.  $\frac{y_{n+1,i} - y_{n+1,i}^*}{\mathbf{sc}_i} = \frac{y_{n+1,i} - y_{n+1,i}^*}{a_{\text{tol}} + |y| \text{rtol}} = 1$

Instead of Eq. (17.2.10) of Press, Teukolsky, Vetterling, Flannery (2007) [36], i.e. Eq. 147, put in a safety factor  $S$ ,  $S < 1$  ( $S$  smaller than 1 by a few percent)

$$(148) \quad h_{n+1} = S h_n \left( \frac{1}{\mathbf{err}_n} \right)^{1/5} \equiv S h_n \left( \frac{1}{\mathbf{err}_n} \right)^{\frac{1}{q+1}}$$

25.3.4. *Step-Control Stability.* cf. pp. 24, IV.II Stiff Problems - One-Step Methods, Hairer and Wanner (2010) [38]

Consider

$$(149) \quad \mathbf{err}_1 = y_1 - \hat{y}_1 \equiv \mathbf{err}_{n+1} = \mathbf{y}_{n+1} - \mathbf{y}_{n+1}^*$$

For the stiff case, consider

$$(150) \quad d_0 = y_0 - y(x_0) \equiv \mathbf{d}_n = \mathbf{y}_n - \mathbf{y}(x_n)$$

and so error estimator satisfies

$$\mathbf{err}_1 \approx E(hJ)d_0 \text{ with } E(z) = R(z) - \hat{R}(z) \text{ i.e.} \\ \mathbf{err}_{n+1} \approx E(hJ)\mathbf{d}_n$$

Consider the numerical process for  $y' = \lambda y$ ,

$$(151) \quad \begin{aligned} \mathbf{y}_{n+1} &= R(h_n \lambda) \mathbf{y}_n \\ \mathbf{err}_n &= E(h_n \lambda) \mathbf{y}_n \\ h_{n+1} &= h_n \left( \frac{\text{tol}}{|\mathbf{err}_n|} \right)^\alpha \end{aligned}$$

cf. Eq. (2.28) of Hairer and Wanner (2010) [38], where  $\mathbf{err}_n$  is the estimated error,  $\alpha = \frac{1}{\hat{p}+1}$ ,  $\hat{p}$  is order of  $\hat{R}(z)$ .

Let  $\eta_n = \log |y_n|$ ,  $\chi_n = \log h_n$ .

So Eq. (2.28) of Hairer and Wanner (2010) [38], i.e. Eq. 151, becomes

$$(152) \quad \begin{aligned} \eta_{n+1} &= \log R + \eta_n = \log |R(\exp(\chi_n) \lambda)| + \eta_n \\ \chi_{n+1} &= \chi_n + \alpha (\log \text{tol} - \log |\mathbf{err}_n|) = \alpha (\gamma - \log |E(\exp(\chi_n))| - \eta_n) + \chi_n \end{aligned}$$

where  $\gamma$  is a constant.

The fixed pt.  $(\eta, \chi)$  satisfies

$$(153) \quad |R(\exp(\chi_n) \lambda)| = 1$$

cf. (2.31) of Hairer and Wanner (2010) [38].

25.3.5. *PI Step Size Control.* cf. pp. 28-29, Hairer and Wanner (2010) [38]

From the wikipedia page for a "PID controller", consider a control function

$$u(t) = K_p e(t) + K_i \int_0^t e(\tau) d\tau + K_d \frac{de(t)}{dt}$$

where  $K_p, K_i, K_d \geq 0$  with

$e(t) \equiv$  error value as difference between desired set point  $\text{SP} \equiv r(t)$ , and

$y(t) \equiv PV =$  measured process variable,  $e(t) = r(t) - y(t)$

Then

$$\dot{u}(t) = K_p \dot{e} + K_i e(t) + K_d \ddot{e}$$

Or, from the point of view of Hairer and Wanner (2010) [38], consider  $\theta(t) \equiv$  departure at time  $t$  of quantity to be controlled from its normal value (i.e.  $e(t)$ ),

Let  $C(t) \equiv$  effect of control,  $-m\theta(t) \equiv$  self-regulating effect such as "vessel in a constant temperature bath."

Suppose

$$\dot{\theta}(t) = C(t) - m\theta(t)$$

Assume (more realistically) that our system has some time lag.

So choose  $C(t)$  s.t. (i.e. replace as our assumption)

$$-\dot{C}(t) = n_1 \theta(t - T) + n_2 \dot{\theta}(t - T) + n_3 \ddot{\theta}(t - T)$$

cf. Eq. (2.40") of Hairer and Wanner (2010) [38].

First term  $n_1 \theta(t - T)$  represents "Integral feedback" (I),  $K_i \int_0^t e(\tau) d\tau$ ,  $K_i e$ ,

2nd. term  $n_2 \dot{\theta}(t - T)$  is "proportional feedback" (P),  $K_p e(t)$ ,  $K_p \dot{e}$

last term  $n_3 \ddot{\theta}(t - T)$  is "Derivative feedback" (D),  $K_d \frac{de(t)}{dt}$ ,  $K_d \ddot{e}$

P term increases constant  $m$ , thus adding extra friction.

So analogously, recall Eq. 152

$$\chi_n = \log h_n$$

$$\chi_{n+1} = \chi_n + (-\alpha)(\log |\mathbf{err}| - \log \text{tol})$$

Let  $\theta(t) \equiv e(t) \Leftrightarrow \log |\mathbf{err}| - \log \text{tol}$  (the "deviation")

Let  $C(t) \equiv u(t) \Leftrightarrow \log h_n \equiv \chi_n$  (control variable)

Then

$$h_{n+1} = h_n \left( \frac{\text{tol}}{|\mathbf{err}_n|} \right)^{n_1}$$

corresponds to

$$\chi_{n+1} - \chi_n = -n_1 (\log |\mathbf{err}| - \log \text{tol})$$

which corresponds to only integral "I" control.

Consider *PI*:

$$\begin{aligned} \chi_{n+1} - \chi_n &= -n_1 (\log |\mathbf{err}_{n+1}| - \log \text{tol}) - n_2 [\log |\mathbf{err}_{n+1}| - \log \text{tol} - (\log |\mathbf{err}_n| - \log \text{tol})] = \\ &= -(n_1 + n_2) (\log |\mathbf{err}_{n+1}| - \log \text{tol}) + n_2 (\log |\mathbf{err}_n| - \log \text{tol}) \end{aligned}$$

$$(154) \quad \implies h_{n+1} = h_n \left( \frac{\text{tol}}{|\mathbf{err}_{n+1}|} \right)^\alpha \left( \frac{|\mathbf{err}_n|}{\text{tol}} \right)^\beta$$

where  $\alpha = n_1 + n_2$ ,  $\beta = n_2$ , i.e.  $\alpha = n_1 + \beta$ .

25.3.6. *Suggestions for values for parameters in Automatic Step-Size Control.*

For **Atol**, **Rtol**, tolerances, Hairer and Wanner (2010) [38], pp. 30 of IV. Stiff Problems - One-Step Methods, for DOPRI5 with  $\beta = 0.13$ , and  $\beta = 0$  (undamped step size control), an extra-large tolerance was chosen to make the difference between the two  $\beta$  values clearly visible, Atol = Rtol=  $8 \cdot 10^{-2}$ . From Press, Teukolsky, Vetterling, Flannery (2007)[36], tldr: first choice is to set **atol** = **rtol** =  $\epsilon = 10^{-6}$ ; Press, Teukolsky, Vetterling, Flannery (2007)[36] says that you may be dealing with a set of equations differing enormously in magnitude; you’d want to use fractional errors, **atol** = 0, **rtol** =  $\epsilon$ , where  $\epsilon$  is the number like  $10^{-6}$  or ”whatever.” If you have oscillatory functions that pass through 0 but are bounded by some max values, you’d want to set **atol** = **rtol** =  $\epsilon = 10^{-6}$ . This choice is safest in general, and should usually be your first choice.

For the *safety factor*, labeled  $S$  or **safe** in pp. 914 of Press, Teukolsky, Vetterling, Flannery (2007)[36], labeled as fac on pp. 168 of Hairer, Nørsett, and Wanner (1993) [37], Press, Teukolsky, Vetterling, Flannery (2007)[36]: safety factor  $S$  should be a few percent smaller than unity (1), Hairer, Nørsett, and Wanner (1993)[37]: fac = 0.8, 0.9,  $(0.25)^{1/(q+1)}$ ,  $(0.38)^{1/(q+1)}$  so error will be acceptable the next time with high probability; further  $h$  is not allowed to increase nor to decrease too fast.

For  $\alpha, \beta$ , Hairer and Wanner (2010) [38], quoted a good choice, found by Gustafsson (1991), after studying small  $h$  and studying a characteristic equation, found

$$(155) \qquad \qquad \qquad \alpha \approx 0.7/p, \quad \beta \approx 0.4/p$$

for  $p$  = exponent of  $h$  of the leading term in the error estimator, cf. (2.48) of Hairer and Wanner (2010) [38]. From pp. 915, 17.2.1 ”PI Stepsize Control” of Press, Teukolsky, Vetterling, Flannery (2007) [36], typically  $\alpha$  and  $\beta$  should be scaled as  $1/k$ , where  $k$  is the exponent of  $h$  in **err**  $\equiv \sqrt{\frac{1}{N} \sum_{i=0}^{N-1} \left( \frac{\Delta_i}{\text{scale}_i} \right)^2}$  ( $k = 5$  for a 5th-order method). Setting  $\alpha = \frac{1}{k}$ ,  $\beta = 0$  recovers the classical controller, Eqn. (17.2.12) of Press, Teukolsky, Vetterling, Flannery (2007) [36] or Eqn. 148,  $h_{n+1} = Sh_n \left( \frac{1}{\text{err}_n} \right)^{1/5}$ . Nonzero  $\beta$  improves stability but loses some efficiency for ”easy” parts of the solution. A good compromise is

$$(156) \qquad \qquad \qquad \beta \approx \frac{0.4}{k}, \quad \alpha \approx \frac{0.7}{k} = \frac{1}{k} - 0.75\beta$$

Be wary, however, of whether this applies to DOPR853, an 8th order Runge-Kutta method. For instance, take a look at the **code from Hairer and Wanner (1994)**, in particular **dopri5.c**:  $\beta$ ’s suggested default value is 0.04 and the formula for  $\alpha$  is  $\alpha = 0.2 - \beta * 0.75$ . The formula corresponds to what was said in Eq. 156. Likewise, for Press, Teukolsky, Vetterling, Flannery (2007) [36], in particular **stepperdopr5.h** found **here**, the default value of  $\alpha$  is also given by Eq. 156.

However, for the DOPR853 implementation, Numerical Recipes suggests for the default value of  $\alpha$  to be given by

$$(157) \qquad \qquad \qquad \frac{1.0}{8.0} - \beta * 0.2 = \frac{1.0}{k} - \beta * 0.2$$

where I speculate that 8 in the denominator is from the order of the Runge-Kutta method in consideration (i.e. DOPR853). Likewise, from Hairer and Wanner’s code, in particular for **dopr853.c**,  $\alpha$  or **expo1** in Hairer and Wanner’s code is given by Eq. 157 as well. The suggested default value is also 0.0 for  $\beta$  and is checked to be between 0.0 and 0.2.

25.3.7. *Embedded Formulas of Order 5.* From Table 5.2. on pp. 178 of Hairer, Nørsett, and Wanner (1993) [37],

$$(158) \qquad \qquad \qquad \begin{array}{c|ccccccc} c_i & a_{ij} \\ \hline 0 & & & & & & & \\ \frac{1}{5} & \frac{1}{5} & & & & & & \\ \frac{3}{10} & \frac{3}{40} & \frac{9}{40} & & & & & \\ \frac{4}{5} & \frac{44}{45} & \frac{-56}{15} & \frac{32}{9} & & & & \\ \frac{8}{9} & \frac{19372}{6561} & \frac{-25360}{2187} & \frac{64448}{6561} & \frac{-212}{729} & & & \\ 1 & \frac{9017}{3168} & \frac{-355}{33} & \frac{46732}{5247} & \frac{49}{176} & \frac{-5103}{18656} & & \\ 1 & \frac{35}{384} & 0 & \frac{500}{1113} & \frac{125}{192} & \frac{-2187}{6784} & \frac{11}{84} & \\ y_n, b_i & \frac{35}{384} & 0 & \frac{500}{1113} & \frac{125}{192} & \frac{-2187}{6784} & \frac{11}{84} & 0 \\ y_n^*, b_i^* & \frac{5179}{57600} & 0 & \frac{7571}{16695} & \frac{393}{640} & \frac{-92097}{339200} & \frac{187}{2100} & \frac{1}{40} \end{array}$$

25.4. **Dense Output, Discontinuities, Derivatives, Interpolation.** cf. II.6 ”Dense Output, Discontinuities, Derivatives”, pp. 188 of Hairer, Nørsett, and Wanner (1993) [37].

Dense output formulas - Runge Kutta methods which provide in addition to numerical result  $y_1 \equiv \mathbf{y}_{n+1}$ , numerical approximation to  $y(x_0 + \theta h) \equiv \mathbf{y}(x_n + \theta h)$ ,  $\forall 0 \leq \theta \leq 1$ .

Start from  $s$ -stage Runge-Kutta method and add  $s^* - s$  new stages, and consider

$$(159) \qquad \qquad \qquad u(\theta) = y_0 + h \sum_{i=1}^{s^*} b_i(\theta) k_i \equiv u(\theta) = \mathbf{y}_n + h \sum_{i=1}^{s^*} b_i(\theta) \mathbf{k}_i$$

cf. Eq. (6.1) of Hairer, Nørsett, and Wanner (1993) [37], where

$$(160) \qquad \qquad \qquad k_i = f(x_0 + c_i h, y_0 + h \sum_{j=1}^{i-1} a_{ij} k_j) = \mathbf{k}_i = \mathbf{f}(x_n + c_i h, \mathbf{y}_n + h \sum_{j=1}^{i-1} a_{ij} \mathbf{k}_j) \quad i = 1, \dots s^*$$

cf. Eq. (6.2)of Hairer, Nørsett, and Wanner (1993) [37]. where  $b_i(\theta)$  are polynomials determined s.t.

$$(161) \qquad \qquad \qquad u(\theta) - y(x_0 + \theta h) = \mathcal{O}(h^{p^*+1})$$

usually  $s^* \geq s + 1$  since we include (at least) the 1st. function evaluation of the subsequent step  $k_{s+1} = hf(x_0 + h, y_1)$  in formula with  $a_{s+1,j} = b_j, \quad \forall j$ .

25.4.1. *Interpolation.* From Sec. 3, pp. 143, ”Interpolation” by Shampine (1986) [39], in addition to the usual step from  $(x_n, \mathbf{y}_n)$  of length  $h$ , we independently advance the integration a step of length  $h^* = \sigma h$  by another formula:

**Theorem 1** (6.1, Hairer, Nørsett, and Wanner (1993) [37]). *Error of approximation  $\mathbf{u}(\theta) = \mathbf{y}_n + h \sum_{i=1}^{s^*} b_i(\theta) \mathbf{k}_i$  is of order  $p^*$  (i.e. local error satisfies Eq. (6.3) of Hairer, Nørsett, and Wanner (1993) [37], i.e.  $\mathbf{u}(\theta) - \mathbf{y}(x_0 + \theta h) = \mathcal{O}(h^{p^*+1})$ ) iff*

$$(162) \qquad \qquad \qquad \sum_{j=1}^{s^*} b_j(\theta) \Phi_j(t) = \frac{\theta^{\varrho(t)}}{\gamma(t)} \text{ for } \varrho(t) \leq p^*$$

with  $\Phi_j(t), \varrho(t), \gamma(t)$  given in Sec. II.2 of Hairer, Nørsett, and Wanner (1993) [37].



*Proof.*  $q$ th derivative (with respect to  $h$ ) of numerical approximation is given by (2.14) with  $b_j$  replaced by  $b_j(\theta)$ ; that of exact solution  $\mathbf{y}(x_0 + \theta h)$  is  $\theta^q y^{(q)}(x_0)$ .

Recall Thm. 2.11 of Hairer, Nørsett, and Wanner (1993) [37],

$$(163) \quad y_1^{(q)} \Big|_{h=0} = \sum_{t \in LT_q} \gamma(t) \sum_j b_j \Phi_j(t) F(t)(y_0) = \sum_{t \in T_q} \alpha(t) \gamma(t) \sum_j b_j \Phi_j(t) F(t)(y_0)$$

cf. Eq. (2.14) of Hairer, Nørsett, and Wanner (1993) [37], which then leads to Th. 2.13 of Hairer, Nørsett, and Wanner (1993) [37].  $\square$

**Corollary 1** (6.2 of Hairer, Nørsett, and Wanner (1993) [37]). *Condition (6.4) of Hairer, Nørsett, and Wanner (1993) [37], Eq. 162, implies derivatives of Eq. (6.1) of Hairer, Nørsett, and Wanner (1993) [37],  $\mathbf{u}(\theta) = \mathbf{y}_n + h \sum_{i=1}^{s^*} b_i(\theta) \mathbf{k}_i$ , Eq. 159 approximate derivatives of exact solution as*

$$(164) \quad h^{-k} \mathbf{u}^{(k)}(\theta) = \mathbf{y}^{(k)}(x_0 + \theta h) = \mathcal{O}(h^{p^* - k + 1})$$

*Proof.* Consider Eq. 159,

$$\mathbf{u}(\theta) - \mathbf{y}(x_n + \theta h) = \mathcal{O}(h^{p^* + 1})$$

which is the error of approximation  $\mathbf{u}(\theta)$ , i.e. local error condition.

While Hairer, Nørsett, and Wanner (1993) [37] says on pp. 189 to take the  $q$ th derivative with respect to  $h$ , I think that's wrong. Take the partial derivative with respect to  $\theta$ :

$$\begin{aligned} \frac{\partial \mathbf{u}(\theta)}{\partial \theta} - \frac{\partial \mathbf{y}}{\partial \theta}(x_n + \theta h) &\equiv \mathbf{u}'(\theta) - h \mathbf{y}'(x_n + \theta h) \xrightarrow{\frac{1}{h}} h^{-1} \mathbf{u}'(\theta) - \mathbf{y}'(x_n + \theta h) = \mathcal{O}(h^{p^* + 1 - 1}) \\ \frac{\partial^k \mathbf{u}(\theta)}{\partial \theta^k} - \frac{\partial^k \mathbf{y}}{\partial \theta^k}(x_n + \theta h) &\equiv \frac{\partial^k \mathbf{u}(\theta)}{\partial \theta^k} - h^k \mathbf{y}^{(k)}(x_n + \theta h) \xrightarrow{\frac{1}{h^k}} h^{-k} \frac{\partial^k \mathbf{u}(\theta)}{\partial \theta^k} - \mathbf{y}^{(k)}(x_n + \theta h) = \mathcal{O}(h^{p^* + 1 - k}) \end{aligned}$$

So Eq. (6.5) of Hairer, Nørsett, and Wanner (1993) [37], i.e. Eq. 164, would be equivalent to saying that, for the  $k = 1$  case,

$$\frac{\partial}{\partial \theta} \sum_{j=1}^{s^*} b_j(\theta) \Phi_j(t) = \sum_{j=1}^{s^*} \frac{\partial b_j(\theta)}{\partial \theta} \Phi_j(t) \equiv \sum_{j=1}^{s^*} b'_j(\theta) \Phi_j(t) = \frac{\varrho(t) \theta^{\varrho(t) - 1}}{\gamma(t)} \text{ for } \varrho(t) \leq p^*$$

But this follows by differentiation. The immediate above has a similar form for  $k > 1$ .  $\square$

25.4.2. *Hermite interpolation.* From (Shampine 1985), Shampine had this easier and more efficient way for low order dense output formulas by using HErmite interpolation. This interpolation method makes use of information produced during integration.

We have 2 function values  $y_0, y_1$  ( $\equiv \mathbf{y}_n, \mathbf{y}_{n+1}$ ) and 2 derivatives

$$\mathbf{f}_0 = \mathbf{f}(x_0, \mathbf{y}_0); \mathbf{f}_1 = \mathbf{f}(x_0 + h, \mathbf{y}_1) \equiv \mathbf{f}_n = \mathbf{f}(x_n, \mathbf{y}_n), \mathbf{f}_{n+1} = \mathbf{f}(x_n + h, \mathbf{y}_n), \text{ i.e. } \mathbf{y} \text{ and } \frac{d\mathbf{y}}{dx} = \mathbf{f}$$

at beginning and end of the step.

Do cubic polynomial interpolation:

Consider Eq. (6.7) on pp. 190 of Hairer, Nørsett, and Wanner (1993) [37],

$$(165) \quad \mathbf{u}(\theta) = (1 - \theta) \mathbf{y}_n + \theta \mathbf{y}_{n+1} + \theta(\theta - 1)((1 - 2\theta)(\mathbf{y}_{n+1} - \mathbf{y}_n) + (\theta - 1)h\mathbf{f}_n + \theta h\mathbf{f}_{n+1})$$

This is the same as Eq. (17.2.15) on pp. 916 of Press, Teukolsky, Vetterling, Flannery (2007) [36]:

$$(166) \quad \mathbf{y}(x_n + \theta h) = (1 - \theta) \mathbf{y}_n + \theta \mathbf{y}_{n+1} + \theta(\theta - 1)[(1 - 2\theta)(\mathbf{y}_{n+1} - \mathbf{y}_n) + (\theta - 1)h\mathbf{f}_n + \theta h\mathbf{f}_{n+1}]$$

where  $0 \leq \theta \leq 1$ .

Whenever underlying method is of order  $p \geq 3$ , we thus obtain a continuous Runge-Kutta method, i.e. a value of  $\mathbf{y}$  that's of order 3, i.e. 3rd. order accurate, as can be verified with a Taylor series.

I had tried to verify this, that the "Hermite interpolation is a special case of (6.1)", pp. 190 of Hairer, Nørsett, and Wanner (1993) [37], by plugging in with  $\mathbf{y}_{n+1} = \mathbf{y}_n + h \sum_{i=1}^s b_i \mathbf{k}_i$ ,

$$\mathbf{u}(\theta) = \mathbf{y}_n + \theta h \sum_{i=1}^s b_i \mathbf{k}_i + \theta(\theta - 1) \left[ (1 - 2\theta)(h \sum_{i=1}^s b_i \mathbf{k}_i) + (\theta - 1)h\mathbf{f}_n + \theta h\mathbf{f}_{n+1} \right]$$

Notice that the Hermite interpolation is essentially this in terms of inputs and outputs:

$$(167) \quad \text{numerical values } \mathbf{y}_n, \mathbf{y}_{n+1}, \frac{d\mathbf{y}}{dx}(x_n), \frac{d\mathbf{y}}{dx}(x_n + h), \xrightarrow[\theta]{\text{Hermite Interpolation with parameter } \theta} \mathbf{u}(\theta)$$

25.4.3. *Dense output implementation.* Consider developing an implementation for dense output.

We can generate the desired points along  $x$  (independent variable, which could be dneoted as  $t$ ), from the endpoints  $x_1, x_2$ , and number of intervals  $N$ ,

$$h = \frac{x_2 - x_1}{N}, \quad x_j = x_1 + hj, \quad \forall j = 0, 1, \dots, N$$

Given a method with stepsize PI control,

$$x_n, y_n, \frac{dy}{dx}(x_n) \mapsto x_{n+1}, y_{n+1}, \frac{dy}{dx}(x_{n+1})$$

Suppose we have run this step and obtained

$$x_{n+1}, y_{n+1}, \frac{dy}{dx}(x_{n+1})$$

and we had saved our initial inputs:

$$x_n, y_n, \frac{dy}{dx}(x_n)$$

Give a  $x$  point in the "dense output",  $x_j$ , we can do hermite interpolation only if

$$x_n \leq x_j \leq x_{n+1}$$

Since  $x_1 \leq x_j$ , let's focus in on this case:  $x_{n+1} < x_j$ . If that's the case, run the step again so that

$$x_j \leq x_{n+1}$$

or until it is so.

Consider checking first that  $x_j \leq x_{n+1}$ . If so, interpolate, and then consider  $x_{j+1}$ . Otherwise, then  $x_j > x_{n+1}$ , so run the step so that  $x_n$  is set to  $x_{n+1}$ , but  $x_{n+1}$  is the new value. Then check again. So we'll need to keep track of  $x_n, x_{n+1}$  for  $x_j$  to be in the interval  $[x_n, x_{n+1}]$

25.5. **Interpreting Runge-Kutta method implementations in Numerical Recipes, and other software.** cf. Press, Teukolsky, Vetterling, Flannery (2007) [36].

25.5.1. *StepperBase*. cf. pp. 903-904, 17.0.2 The Odeint Object, Press, Teukolsky, Vetterling, Flannery (2007) [36]

Inputs of the class StepperBase are the following:

$\mathbf{y}_0 = \mathbf{y}_0(x_0) \in \mathbb{R}^n$   
 $\frac{d\mathbf{y}}{dx}(x_0) \in \mathbb{R}^n$   
 $x_0 = x(0) \in \mathbb{R} \equiv$  starting value of independent variable  $x$   
 $\epsilon_a \equiv$  atol  
 $\epsilon_r \equiv$  rtol

For pp. 918, Ch. 17 Integration of Ordinary Differential Equations, Press, Teukolsky, Vetterling, Flannery (2007) [36], for the algorithm routine `dy`, it computes the 6 steps to compute  $\mathbf{k}_i$ ,  $i = 1, \dots, 6$  and the seventh FSAL step (which I believe computes the next  $y(x_{n+1})$  value anyway), and computes  $y_{n+1}$  and error  $\Delta$ .

So for `stepperdopr5.h` code on pp. 918,

```
template <class D>
void StepperDopr<D>::dy(const Doub h, D &derivs), then
    c2, ... c5  $\equiv c_2, \dots, c_5$  and  $c_6 = c_7 = 1$  are not explicitly defined in the code.
a21, a31, a32, ... a75, a76  $\equiv a_{21}, a_{31}, a_{32} \dots a_{75}, a_{76}$  coefficients.
e1, e3, ... e6, e7  $\equiv b_1^*, b_3^*, \dots, b_6^*, b_7^*$ , where  $b_2^* = 0$  isn't explicitly defined in the code.
```

$f \equiv$  `derivs` (in the code)

```
    k1 = f(xn, yn) =  $\frac{dy}{dx}(x_n) \equiv$  dydx[i]
ytemp[i] = yn + ha21k1  $\equiv y_n + ha_{21} \frac{dy}{dx}(x_n)$ 
k2 = f(xn + c2h, ytemp[i])
ytemp[i] = yn + h(a31  $\frac{dy}{dx}(x_n)$  + a32k2)
k3 = f(xn + c3h, ytemp[i])
:
k6 = f(xn + c6h, yn + h(a61  $\frac{dy}{dx}(x_n)$  + a62k2 + ... + a65k5))
where c6 = 1.
```

In the code, `yout` corresponds to the following in our notation:

`yout`  $\equiv y(x_{n+1})$

$$y(x_{n+1}) = y(x_n) + h(a_{71} \frac{dy}{dx}(x_n) + a_{72}k_2 + \dots + a_{76}k_6) = y(x_n) + h(b_1a_{71} + \dots + b_7a_{76})$$

Thus, essentially, the inputs and outputs of this function, `StepperDopr<D>::dy(...)` is this:

$$(168) \quad \begin{array}{c} \text{constants } c_2, \dots, c_2, a_{ij}; b_i^* \\ h, f: \mathbb{R}, \mathbb{R}^N \rightarrow \mathbb{R}^N \\ \text{initial values } \frac{dy}{dx}(x_n), x_n, \mathbf{y}_n \end{array} \xrightarrow[\text{Runge-Kutta method}]{\text{dy} \equiv} \mathbf{y}_{n+1}, \mathbf{f}(x_n + h, \mathbf{y}_{n+1}) \equiv \frac{d\mathbf{y}}{dx}(x_n + h); \Delta \equiv \mathbf{y}_{n+1} - \mathbf{y}_{n+1}^*$$

25.5.2. *prepare dense, dense output*. Consider the coefficients shown in the code for `StepperDopr5<D>::prepare_dense()` on pp. 918 of Press, Teukolsky, Vetterling, Flannery (2007) [36]:

(169)

$$\begin{aligned} d_1 &= \frac{-12715105075.0}{11282082432.0} \\ d_2 &= 0 \\ d_3 &= \frac{87487479700.0}{32700410799.0} \\ d_4 &= \frac{-10690763975.0}{1880347072.0} \\ d_5 &= \frac{701980252875.0}{199316789632.0} \\ d_6 &= \frac{-1453857185.0}{822651844.0} \\ d_7 &= \frac{69997945.0}{29380423.0} \end{aligned}$$

Someone had wondered how these dense output coefficients in Eq. 169 were derived: <https://math.stackexchange.com/questions/2947231/how-can-i-derive-the-dense-output-of-ode45>

Furthermore, for the function `prepare_dense()` it finally utilizes these class data members (EY: which may not be necessary as class data members?) `rcont1`, `rcont2`, ...; I present the code and then use my notation:

```
    rcont1[i] = y[i];  $\rightarrow$  rcont1 = y
rcont2[i] = ydiff;  $\rightarrow$  rcont2 =  $\mathbf{y}_{\text{out}} - \mathbf{y}$ 
Doub bsp1 = h * dydx[i] - ydiff;  $\rightarrow$  rcont3 =  $h \frac{dy}{dx} - (\mathbf{y}_{\text{out}} - \mathbf{y})$ 
rcont4[i] = ydiff-h*dydxnew[i] - bsp1  $\rightarrow$  rcont4 =  $(\mathbf{y}_{\text{out}} - \mathbf{y}) - h \frac{d\mathbf{y}}{dx}(x_n + h) - h \frac{d\mathbf{y}}{dx}(x_n) + (\mathbf{y}_{\text{out}} - \mathbf{y}) = 2(\mathbf{y}_{n+1} - \mathbf{y}) -$ 
 $h \left( \frac{d\mathbf{y}}{dx}(x_n + h) + \frac{d\mathbf{y}}{dx}(x_n) \right)$ 
rcont5  $\rightarrow h \left[ d_1 \frac{d\mathbf{y}}{dx}(x_n) + d_3 \mathbf{k}_3 + d_4 \mathbf{k}_4 + d_5 \mathbf{k}_5 + d_6 \mathbf{k}_6 + d_7 \frac{d\mathbf{y}}{dx}(x_n + h) \right]$ 
```

From Shampine (1986) [39], for what was denoted as  $c_j^*$  - result at  $x_n + \frac{1}{2}h$ , order 4, the  $c_j^*$  coefficients (Shampine numbers from  $j = 0, 1, \dots$  while we will number from  $j = 1, 2, \dots$ )

$$\begin{aligned} &c_j^* \text{ result at } x_n + \frac{1}{2}h \\ &\frac{6025192743}{30085553152} \\ &0 \\ &\frac{51252292925}{65400821598} \\ &\frac{-2691868925}{45128329728} \\ &\frac{187940372067}{1594534317056} \\ &\frac{-1776094331}{19743644256} \\ &\frac{11237099}{235043384} \end{aligned}$$

Take a look on pp. 918-919 of Press, Teukolsky, Vetterling, Flannery (2007) [36] for the implementation of `void StepperDopr5<D>::prepare_dense(const Doub h, D& derivs)` and `Doub StepperDopr5<D>::dense_out(const Int i, const Doub x, const Doub h);`

the return value of `StepperDopr5<D>::dense_out(...)` is

$$\mathbf{y}_n + \theta(\mathbf{y}_{n+1} - \mathbf{y}_n + (1 - \theta)(h \frac{d\mathbf{y}}{dx}(x_n) - (\mathbf{y}_{n+1} - \mathbf{y}_n) + \theta(2(\mathbf{y}_{n+1} - \mathbf{y}) - h \left[ \frac{d\mathbf{y}}{dx}(x_n + h) + \frac{d\mathbf{y}}{dx}(x_n) \right] ))) + \theta^2(1 - \theta)^2(h \sum_{i=1}^{s+1} d_i \mathbf{k}_i)$$

Take a look at all the terms above except for the last one. It is essentially the *Hermite interpolation* found in Eq. [165](#):

$$(1 - \theta)\mathbf{y}_n + \theta\mathbf{y}_{n+1} + \theta(\theta - 1) \left[ (1 - 2\theta)(\mathbf{y}_{n+1} - \mathbf{y}) + (\theta - 1)h \frac{d\mathbf{y}}{dx}(x_n + h) + \theta h \frac{d\mathbf{y}}{dx}(x_n) \right] = \mathbf{u}(\theta)$$

To refactor this, consider the inputs and outputs needed from Eq. [167](#).

25.5.3. *Controller success is PI Control Step Size.* Consider the class member function `bool StepperDopr5<D>::Controller::success(const Doub err, Doub& h)` on pp. 919 of Press, Teukolsky, Vetterling, Flannery (2007). Rewrite it in pseudocode:

```
    if err == 0
scale = maxscale
if 0 < err ≤ 1.0
scale = safe · err−α · eroldβ
if scale < minscale
scale = minscale
if scale > maxscale
scale = maxscale
```

```
    if (reject)
hnext = h · min(scale, 1.0)
else
hnext = h · scale
erold = max(err, 10−4)
reject = false
return true
```

```
    else (if err > 0)
scale = max(safe · err−α, minscale)
h ↦ h · scale
reject = true
return false
```

Let’s see where `success(...)` method gets used. Note that essentially, `bool StepperDopr<D>::Controller::success(const Doub err, Doub& h) ↦ reject, erold, hnext`, i.e. this method affects these 3 data members of the `struct Controller`. Indeed, the `struct Controller` has these data members, `hnext, erold, reject`.

This method gets called in `con.success(err, h)` and appears to only be in this class method - `void StepperDopr5<D>::step(const Doub htry, D& derivs)`. And so if we look at pp. 917 for the code right above, then `Doub err = error()`, i.e. it’s the scaled error estimate, given by Eq. [145](#), that gets fed into the class method `success`.

25.5.4. *step, the actual stepper.* Recall some of the code for `StepperDopr5<D>::step(...)` for `stepperdopr5.h` on pp. 917, Ch. 17. ”Integration of Ordinary Differential Equations” of Press, Teukolsky, Vetterling, Flannery (2007):

```
template <class D>
void StepperDopr5<D>::step(const Doub htry, D& derivs)
{
    Doub h = htry;
```

```
    for (;;)
    {
        dy(h, derivs);
        Doub err=error();
        if (con.success(err, h)) break;
        ...
    }
    ...
}
```

where `success`, recall, passes by reference the local variable `h`, which I personally don’t advise. But let’s figure out what `success` is doing. It mutates (again, a programming no-no) `hnext` if the error is less than or equal to 1.0 and mutates our reference to `h` otherwise.

25.5.5. *Interpreting code for DOPR853 methods.* Consider these implementations of the DOPR853 method provided for [Numerical Recipes](#) and from [Hairer](#).

In both implementations, the only coefficients involved with calculating  $\mathbf{y}_{n+1}$  and  $\frac{d\mathbf{y}}{dx}(x_{n+1})$  are

$$a_{21}, a_{31}, a_{32}, \dots, a_{12,1}, a_{12,2} \dots a_{21,11}$$

and

$$b_1, \dots, b_6, b_7, \dots b_{12}$$

We *could* have, or had the possibiltiy of, copying the values of  $b_i$  into another ”row” of  $a$  coefficients as used in the ”Embedded Formulas” scheme of Hairer, Nørsett, and Wanner (1993) [\[37\]](#), but we will not; so we’ll follow Eq. [123](#)

$$\mathbf{k}_j = \mathbf{f}(x_n + c_j h, \mathbf{y}_n + h \sum_{i=1}^{j-1} a_{j,i} \mathbf{k}_i) \forall j = 1, \dots s$$

(171)

$$\mathbf{y}(x_{n+1}) = \mathbf{y}_n + h \sum_{i=1}^s b_i \mathbf{k}_i$$

with  $s = 12$  for DOPR853. Likewise, for the `c` coefficients that are used in the Runge-Kutta method in calculating the value of the independent or ”x” (or ”t”) value, while both codes provide `c` coefficients indexed from 2 to 16, we’ll only use the coefficients indexed from 2 to 11, and add a  $c_{12} = 1$ .

Consider  $\mathbf{yerr} \equiv \mathbf{y}_{err}$

$$\mathbf{y}_{err} = \left( \frac{\mathbf{y}_{out} - \mathbf{y}_n}{h} \right) - b_{hh1} \frac{d\mathbf{y}}{dx} - b_{hh2} \mathbf{k}_9 - b_{hh3} \mathbf{k}_{12}$$

$$\mathbf{yerr2} \equiv \mathbf{y}_{err,2}$$

$$\mathbf{y}_{err2} = e_1 \mathbf{k}_1 + e_6 \mathbf{k}_6 + e_7 \mathbf{k}_7 + \dots + e_{11} \mathbf{k}_{11} + e_{12} \mathbf{k}_{12}$$

$$\mathbf{err2} = \sum_{i=1}^n \left( \frac{y_{err,i}}{sk_i} \right)^2$$
$$\mathbf{err} = \sum_{i=1}^n \left( \frac{y_{err2,i}}{sk_i} \right)^2$$

### 25.6. References for Runge-Kutta method.

- <https://www.unige.ch/~hairer/software.html>

## 26. SYSTEMS OF ORDINARY DIFFERENTIAL EQUATIONS ODES

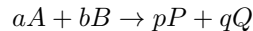
EY : 20160228 There’s a lack of fully general and useful examples of solving a system of ordinary differential equations (ODEs) in `scipy` if one does a search on Google. There is the [scipy cookbook](#), which has toy examples for the [Coupled spring-mass system](#), [Korteweg de Vries equation](#), [Matplotlib: lotka volterra tutorial](#), [Modeling a Zombie Apocalypse](#), [Theoretical ecology: Hastings and Powell](#). Then, there’s Kevin Dunn’s material for Process Model Formulation and Solution, which is *excellent* for its layout, thoroughness, side-by-side comparison between Matlab and Python (`scipy`), and *real-world* examples, though related to *chemical engineering* (but which is an interesting and useful subject in its own right; so we’d might as well learn some chemical engineering).

There was some [MIT OCW material on CSTR](#), but it only solved for the steady-state solution.

**26.1. Continuously Stirred Tank Reactor (CSTR).** Consider the *continuously stirred tank reactor* (CSTR) or vat or backmix reactor.

The main assumption is of perfect mixing - perfect mixing throughout, with each locality’s mixing the same as another locality, so reaction rate is equal everywhere throughout a volume  $V$ .

From [wikipedia’s Reaction Rate](#) article, for a general chemical reaction of the form



the reaction rate  $r$  is given by

$$r = \frac{-1}{a} \frac{d[A]}{dt} = -\frac{1}{b} \frac{d[B]}{dt} = \frac{1}{P} \frac{d[P]}{dt} = \frac{1}{q} \frac{d[Q]}{dt}$$

where  $[X] \equiv$  concentration of substance  $X$ .

26.1.1. *Second-Order Reaction Rate.* A second order reaction rate would be of the form

$$r = kC_A^2$$

and if  $[C_A] = \text{mol}/L$ , then the units for  $k$  are  $[k] = \frac{1/s}{\text{mol}/L} = \frac{L}{\text{mol} \cdot s}$ .

Let’s count (or account) for number of particles  $N$ . In general, for concentration  $n := \frac{N}{V}$ , and  $N = N(t)$ , dependent upon time, then

$$N = \int_V n \text{vol}^n$$

$$\frac{dN}{dt} \equiv \dot{N} = \int_V \frac{\partial n}{\partial t} \text{vol}^n + \int_V \text{div}(nu) \text{vol}^n = \int_V \frac{\partial n}{\partial t} \text{vol}^n + \int_{\partial V} nu^j dS_j$$

One can interpretation the terms of the last equation immediately above as [accumulation] = [generation] + ([in] − [out]), respectively.

Assuming perfect mixing, so the reaction rate equal everywhere throughout volume  $V$ , so  $\frac{\partial n}{\partial t}$  the same throughout  $V$ , and that  $\partial V = S_{\text{in}} \amalg S_{\text{out}}$ , then

$$\frac{dN}{dt} = \nu k C^2 V + \int_{S_{\text{in}}} nu^j dS_j + \int_{S_{\text{out}}} nu^j dS_j$$

Now for  $F(t)$  being the inlet flow of volume per second or volume per minute, then

$$\int_{S_{\text{in}}} nu^j dS_j = F(t) C_{A,\text{in}}$$

Supposing the flow rate out, *in terms of volume*, is the same as flow rate in,

$$\int_{S_{\text{out}}} nu^j dS_j = -F(t) C_A$$

i.e., after perfect mixing, obtaining a concentration  $n = C_A$ , then multiply it by the volume per second or volume per minute flowing out,  $-F(t)$  to obtain the number of particles rushing out.

I think the flow rate out being the same as flow rate in, *in terms of volume*, is what Kevin Dunn meant in the [first example problem for ODE integration](#).

Thus, after dividing by fixed  $V$  (the volume of the tank),

$$\xrightarrow{\frac{1}{V}} \frac{dn_A}{dt} \equiv \frac{dC_A}{dt} = \frac{F(t)C_{A,\text{in}}}{V} - \frac{F(t)C_A}{V} + \nu_A k C_A^2$$

For the “right-hand side” (RHS) of an ODE with constant coefficients, and a RHS that doesn’t depend on time, then solving this problem is pretty easy. In the language of flows on a manifold, letting curve, or our solution, be  $\gamma : \mathbb{R} \rightarrow M$ , so  $\dot{\gamma} \in \mathfrak{X}(M)$ , a vector field on smooth manifold  $M$ , then

$$\frac{d\gamma}{dt} = f(\gamma, t) = f(\gamma) \implies \frac{d\gamma}{f(\gamma)} = dt$$

A copy of Dunn’s implementation is in `CSTRconstant.py` and `CSTRdynamic.py`, copied from verbatim.

EY : 20160229 Things I need to understand further are the `set_integrator` method of `scipy.integrate.ode`, with option `’vode’` and method `’bdf’` (i.e. `method=’bdf’`) and what a [Stiff Equation](#) is.

## 27. COMPUTATIONAL FLUID DYNAMICS (CFD) AND NAVIER-STOKES EQUATION

Consider Liquid Bipropellant Spray Combustion Modeling (liquid bipropellant in that we are considering a liquid fuel, that’s injected in, and also a liquid oxidizer, injected in), as outlined by Preclik, Knab, Estublier, and Wannerberg (Preclik, et. al.) in their article “Simulation and Analysis of Thrust Chamber Flowfields: Storable Propellant Rockets” (2004) [\[23\]](#).

Consider first gas-phase flow modeling. They treat the gaseous phase in 2-dimensions, specifically in an axisymmetric coordinate axes setup, and solving, on it, the Favre-averaged Navier-Stokes equations, with species continuity, and so-called  $k - \epsilon$  turbulence equations. Let’s examine first Navier-Stokes equations.

**27.1. Navier-Stokes equations.** Cuong Nguyen’s notes (November 5, 2005) on [Turbulence Modeling](#) has an incomplete review of what RANDS is, what Favre-averaged Navier-Stokes equations is, and what  $k - \omega$ ,  $k - \epsilon$  are.

27.2. Spacetime  $\iff$  Grid.

27.2.1. *Spatial (smooth)  $N$  manifold  $\iff$  Grid.*

$$\begin{array}{ccc} \mathbb{R}^3 & \xrightarrow{\text{discretization}} & \mathbb{Z}^3 \\ (l_x, l_y, l_z) \in \mathbb{R} \times \mathbb{R} \times \mathbb{R} & \xleftarrow{l_i = N_i h_i} & (N_x, N_y, N_z) \in \mathbb{Z}^+ \times \mathbb{Z}^+ \times \mathbb{Z}^+ \\ & & \downarrow \\ & & (h_x, h_y, h_z) \in (\mathbb{R}^+)^3 \\ (x, y, z) \in \mathbb{R}^3 & \xleftarrow{x^i = i_i h_i} & (i_x, i_y, i_z) \in \{0 \dots N_x - 1\} \times \{0 \dots N_y - 1\} \times \{0 \dots N_z - 1\} \end{array}$$

There is a *flatten* functor that is necessitated by either the contiguous architecture of memory addresses on memory of the device GPU or by software constraints (CUDA C/C++ 7.5 Toolkit doesn’t take multidimensional arrays).

$$\begin{array}{ccc} \mathbb{Z}^3 & \xrightarrow{\text{flatten}} & \mathbb{Z} \\ (N_x, N_y, N_z) \in \mathbb{Z}^+ \times \mathbb{Z}^+ \times \mathbb{Z}^+ & \xrightarrow{\text{flatten}} & N_x N_y N_z \in \mathbb{Z} \\ (i_x, i_y, i_z) \in \{0 \dots N_x - 1\} \times \{0 \dots N_y - 1\} \times \{0 \dots N_z - 1\} & \xrightarrow{\text{flatten}} & i_x + i_y N_x + i_z N_x N_y \equiv i \in \{0 \dots N_x N_y N_z - 1\} \end{array}$$

Note that  $i$  is sometimes denoted as the “global” index in as it directly accesses the memory address on the device (GPU).

Consider  $\rho^\infty(\mathbb{R}^3) \in C^\infty(\mathbb{R}^3)$  and its behavior under the discretization (“discretize”) and flatten functors. Also, treat  $C^\infty(\mathbb{R}^3)$  as the “zero”th order (trivial) vector bundle, endowed with a vector space structure itself (it’s a ring, I believe, and it sits on the spatial manifold  $N = \mathbb{R}^3$ ). Then there is a natural projection  $\pi$  back onto  $N$ .

$$\begin{array}{ccccc}
C^\infty(\mathbb{R}^3) & \ni \rho(x) \in C^\infty(\mathbb{R}^3) & \xrightarrow{\text{discretize}} & \text{Hom}(\mathbb{Z}^3, \mathbb{R}) & \xrightarrow{\text{flatten}} & \text{Hom}(\mathbb{Z}, \mathbb{R}) \\
\pi \downarrow & \pi \downarrow & & \pi \downarrow & & \pi \downarrow \\
\mathbb{R}^3 & \ni x \in \mathbb{R}^3 & \xrightarrow{\text{discretize}} & (i_x, i_y, i_z) & \xrightarrow{\text{flatten}} & i \in \{0 \dots N_x N_y N_z - 1\}
\end{array}$$

Note that we can say that the discretization of  $\rho(\mathbb{R}^3) \in C^\infty(\mathbb{R}^3)$  is a homomorphism  $\text{Hom}$  from  $\mathbb{Z}^3$  to  $\mathbb{R}$  because vector space structure is preserved (and so discretization (discretize) is a functor, along with flatten). As  $C^\infty(N)$  is a vector space (ring) over the field  $\mathbb{R}$  in that it is equipped with commutative (abelian) addition and scalar multiplication by field  $\mathbb{R}$

$$f(x) + g(x) = (f + g)(x) \xrightarrow{\text{discretize}} f(i_x, i_y, i_z) + g(i_x, i_y, i_z) = (f + g)(i_x, i_y, i_z)$$

$$\lambda f(x) \in C^\infty(\mathbb{R}^3) \xrightarrow{\text{discretize}} \lambda f(i_x, i_y, i_z) \in \text{Hom}(\mathbb{Z}^3, \mathbb{R})$$

Now consider the other object we need to consider, the vector bundle, namely the tangent bundle  $TN$  over spatial (smooth) manifold  $N$ . Namely consider the *vector field*, representing the *velocity vector field*  $u$  as a section of the tangent bundle  $T\mathbb{R}^3$  over Euclidean space (as a smooth manifold)  $\mathbb{R}^3$ . Since  $\dim \mathbb{R}^3 = 3$ , then the vector space  $V \equiv T_x N$  “over a fiber” is of dimension 3, i.e.  $\dim T_x N, \forall x \in N$ .

$$u \in \mathfrak{X}(\mathbb{R}^3) = \Gamma(T\mathbb{R}^3)$$

Recall the structure of a vector bundle, defined with a local trivialization  $\varphi$  on an open set  $U \subset N$ :

$$\begin{array}{ccc}
TN & \supset \pi^{-1}(U) & \xrightarrow{\varphi} U \times V \\
\pi \downarrow & \uparrow \pi^{-1} & \swarrow \text{proj} \\
N & \supset U &
\end{array}$$

Then consider its behavior under discretization (discretize functor):

$$\begin{array}{ccc}
u(x) \in \Gamma(TN) & \xrightarrow{\varphi} & (x, u(x)) \\
\pi \downarrow & \swarrow \text{proj} & \\
x \in N & & 
\end{array}
\quad \xrightarrow{\text{discretize}} \quad
\begin{array}{ccc}
u(i_x, i_y, i_z) \in \text{Hom}(\mathbb{Z}^3, \mathbb{R}^3) & \xrightarrow{\varphi} & ((i_x, i_y, i_z), u(i_x, i_y, i_z)) \in \mathbb{Z}^3 \times \text{Hom}(\mathbb{Z}^3, \mathbb{R}^3) \\
\pi \downarrow & \swarrow \text{proj} & \\
(i_x, i_y, i_z) & & 
\end{array}$$

### 27.3. Finite Volume. *Keywords:* Finite Volume, Upwind methods

Ferziger and Peric (2013) [21]. However, here is my generalization.

#### 27.3.1. Upwind method. Consider the so-called “upwind method.” Consider the 1-dimensional case.

Let  $C_i^1 \equiv i$ th cell of dimension 1, for  $i = 0, \dots, N - 1$ . So there are  $N$  total cells.

Cells are centered at  $x_{2i+1} = l \frac{(2i+1)}{2} = l(i + \frac{1}{2})$ .  $l$  is the 1-dimensional size or length of a single cell. Notice that in this case, I am assuming a uniform grid. Note that this can be easily generalized to a grid with different cell sizes for each cell.

For the  $i$ th cell, which is a 1-(cubic) simplex, a line segment,  $C_i^1$ , it has 2 0-(cubic) simplices (faces), which in this 1-dimensional case, it's 2 isolated points:  $\partial C_i^1 = \{C_{i\pm 1}^0\}$ .

The centers of these faces, i.e. the position of these 2 points, at the ends of the line segment, are

$$x_{C_{i\pm 1}^0} = l \left( \frac{2i+1 \pm 1}{2} \right) = l(i + \frac{1}{2} \pm \frac{1}{2}) = \{li, l(i+1)\}$$

I will take the mass conservation equation, in its integral form, as an example here, but this example can be easily generalized to the convection of any other conserved quantity. Define the average mass density  $\bar{\rho}_i$ :

$$(172) \quad \bar{\rho}_i := \frac{1}{l} \int_{C_i^1} \rho \text{vol}^1$$

For the mass conservation equation (in integral form),

$$(173) \quad \int_V \frac{\partial \rho}{\partial t} = - \int_{\partial V} i_{\mathbf{u}} \rho \text{vol}^1$$

where the integral is taken over the volume  $V$ , and over its boundary  $\partial V$  (which is the surface of  $V$ ).

For the left-hand side (LHS) of Eq. 173, rewrite it in terms of  $\bar{\rho}_i$ ,

$$\begin{aligned}
\int_{C_i^1} \frac{\partial \rho}{\partial t}(t, x) \text{vol}^1 &\approx \int_{C_i^1} \frac{\rho(t + \Delta t, x) - \rho(t, x)}{\Delta t} \text{vol}^1 = \\
&= \frac{1}{\Delta t} \left[ \int_{C_i^1} \rho(t + \Delta t, x) \text{vol}^1 - \int_{C_i^1} \rho(t, x) \text{vol}^1 \right] = \\
&= \frac{l}{\Delta t} [\bar{\rho}_i(t + \Delta t) - \bar{\rho}_i(t)]
\end{aligned}$$

Considering the mass flux through the “surface” or through the endpoints of the line segment, that is a cell in the 1-dimensional case,

$$\int_{\partial C_i^1} i_{\mathbf{u}} \rho \text{vol}^1 = \int_{\partial C_i^1} \rho u^i dS_i$$

then the so-called “upwind” scheme is this:

$$\begin{aligned}
\int_{C_{i+1}^0} \rho u^i dS_i &= u^x(x_{C_{i+1}^0}) \int_{C_{i+1}^0} \rho dS_x = \begin{cases} \bar{\rho}_i u^x(x_{C_{i+1}^0}) & \text{if } u^x(x_{C_{i+1}^0}) > 0 \\ \bar{\rho}_{i+1} u^x(x_{C_{i+1}^0}) & \text{if } u^x(x_{C_{i+1}^0}) < 0 \end{cases} \\
\int_{C_{i-1}^0} \rho u^i dS_i &= -u^x(x_{C_{i-1}^0}) \int_{C_{i-1}^0} \rho dS_x = \begin{cases} -\bar{\rho}_{i-1} u^x(x_{C_{i-1}^0}) & \text{if } u^x(x_{C_{i-1}^0}) > 0 \\ -\bar{\rho}_i u^x(x_{C_{i-1}^0}) & \text{if } u^x(x_{C_{i-1}^0}) < 0 \end{cases}
\end{aligned}$$

For the 3-dimensional case, I refer back to my notes on **Computational Physics** in the 3-dim. “Upwind” subsection.

For a rectangular prism (cubic),

for cell  $C_{ijk}^3$ ,  $i = 0 \dots N_x - 1$ ,  $j = 0 \dots N_y - 1$ ,  $k = 0 \dots N_z - 1$ ,  $N_x \cdot N_y \cdot N_z$  total cells.

Cells centered at

$$(x_{2i+1}, y_{2j+1}, z_{2k+1}) = (l^x \frac{(2i+1)}{2}, l^y \frac{(2j+1)}{2}, l^z \frac{(2k+1)}{2}) = \left( \sum_{l=0}^{i-1} l_l^x + \frac{l_i^x}{2}, \sum_{l=0}^{j-1} l_l^y + \frac{l_j^y}{2}, \sum_{l=0}^{k-1} l_l^z + \frac{l_k^z}{2} \right)$$

For the 3-(cubic) simplex,  $C_{ijk}^3$ , it has 6 2-(cubic) simplices (faces). So for  $C_{ijk}^3$ , consider  $\{C_{i\pm 1, jk}^2, C_{ij\pm 1, k}^2, C_{ijk\pm 1}^2\}$ .

The center of these faces, such as for  $C_{i\pm 1, jk}^2$ ,  $x_{C_{i\pm 1, jk}^2}$ , for instance,

$$x_{C_{i\pm 1, jk}^2} = (x_{2i+1\pm 1}, y_{2j+1}, z_{2k+1}) = (l^x \left( \frac{2i+1 \pm 1}{2} \right), l^y \frac{(2j+1)}{2}, l^z \frac{(2k+1)}{2}) = \left( \sum_{l=0}^{\frac{2i-1\pm 1}{2}} l_l^x, \sum_{l=0}^{j-1} l_l^y + \frac{l_j^y}{2}, \sum_{l=0}^{k-1} l_l^z + \frac{l_k^z}{2} \right)$$

We want the flux. So for

$$\bar{\rho}_{ijk} := \frac{1}{l_i^x l_j^y l_k^z} \int_{C_{ijk}^3} \rho \text{vol}^3$$



then the flux through 2-(cubic) simplices (faces),  $\int \rho i_{\mathbf{u}} \text{vol}^3$ ,

$$\begin{aligned} \int_{C_{i+1,jk}^2} \rho i_{\mathbf{u}} \text{vol}^3 &= \begin{cases} l_j^y l_k^z \bar{\rho}_{ijk} u^x(x_{C_{i+1,jk}^2}) & \text{if } u^x(x_{C_{i+1,jk}^2}) > 0 \\ l_j^y l_k^z \bar{\rho}_{i+1,jk} u^x(x_{C_{i+1,jk}^2}) & \text{if } u^x(x_{C_{i+1,jk}^2}) < 0 \end{cases} \\ \int_{C_{i-1,jk}^2} \rho i_{\mathbf{u}} \text{vol}^3 &= \int_{C_{i-1,jk}^2} \rho u^i dS_i = \int_{C_{i-1,jk}^2} \rho u^i \frac{\sqrt{g}}{(3-1)!} \epsilon_{ii_2 i_3} dx^{i_2} \wedge dx^{i_3} = -u^x(x_{C_{i-1,jk}^2}) \int_{C_{i-1,jk}^2} \rho dy dz = \\ &= \begin{cases} -l_j^y l_k^z \bar{\rho}_{i-1,jk} u^x(x_{C_{i-1,jk}^2}) & \text{if } u^x(x_{C_{i-1,jk}^2}) > 0 \\ -l_j^y l_k^z \bar{\rho}_{i,jk} u^x(x_{C_{i-1,jk}^2}) & \text{if } u^x(x_{C_{i-1,jk}^2}) < 0 \end{cases} \end{aligned}$$

and so on.

**27.4. Finite Difference.** Consider mass conservation.

$$\frac{d}{dt} M = \frac{d}{dt} \int_V m \equiv \frac{d}{dt} \int \rho \text{vol}^d = \int \mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}} \rho \text{vol}^d = \int \frac{\partial \rho}{\partial t} \text{vol}^d + \mathbf{d}u_{\mathbf{u}} \rho \text{vol}^d = \int \left( \frac{\partial \rho}{\partial t} + \text{div}(\rho \mathbf{u}) \right) \text{vol}^d = 0$$

From here,

$$\frac{\partial \rho}{\partial t} + \text{div}(\rho \mathbf{u}) = \frac{\partial \rho}{\partial t} + \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^k} (\rho u^k \sqrt{g}) = \frac{\partial \rho}{\partial t} + \rho \text{div} \mathbf{u} + u^k \frac{\partial \rho}{\partial x^k} = 0$$

If  $\text{div} \mathbf{u} = 0$ , this is the *incompressible* case.

If  $\sqrt{g} = 1$ ,  $\frac{\partial \rho}{\partial t} + \frac{\partial}{\partial x^k} (\rho u^k) = 0$

$$\begin{aligned} 27.4.1. \text{ 2 and more species, mass conservation.} \quad & \text{Suppose } M_{\text{tot}} = M_A + M_B, \\ & \rho_A = \frac{M_A}{|\int \text{vol}^d|}, \quad (\text{assume, cf. Landau and Lifshitz,} \\ & \rho_B = \frac{M_B}{|\int \text{vol}^d|} \end{aligned}$$

that species  $A, B$  “don’t talk to each other”, that is, interaction between them is negligible in terms of occupying volume; this is what gases do).

$$\begin{aligned} \dot{M}_{\text{tot}} &\equiv \frac{d}{dt} M_{\text{tot}} = \frac{d}{dt} (M_A + M_B) = \frac{d}{dt} \left( \int \rho_A \text{vol}^d + \int \rho_B \text{vol}^d \right) = \frac{d}{dt} \left( \int (\rho_A + \rho_B) \text{vol}^d \right) = \\ &= \int \left( \frac{\partial \rho_A}{\partial t} + \text{div}(\rho_A \mathbf{u}) + \frac{\partial \rho_B}{\partial t} + \text{div}(\rho_B \mathbf{u}) \right) \text{vol}^d = \int \frac{\partial (\rho_A + \rho_B)}{\partial t} + \text{div}((\rho_A + \rho_B) \mathbf{u}) = 0 \end{aligned}$$

$$\begin{aligned} (174) \quad & U = U(\sigma, V, N) \\ & dU = \tau d\sigma - p dV + \mu dN \end{aligned}$$

$$\begin{aligned} (175) \quad & H = H(\sigma, p, N) := U + pV \\ & dH = dU + p dV + V dp = \tau d\sigma + \mu dN + V dp \end{aligned}$$

$$\begin{aligned} (176) \quad & G = G(\tau, p, N) = U + pV - \tau \sigma = H - \tau \sigma = \mu N \\ & dG = dU + p dV + V dp - \tau d\sigma - \sigma d\tau = \tau d\sigma - p dV + \mu dN + p dV + V dp - \tau d\sigma - \sigma d\tau = \\ & \quad = -\sigma d\tau + V dp + \mu dN = \mu dN + N d\mu \end{aligned}$$

**27.5. Compressible, inviscid (no viscosity) flow i.e. Euler equations.** Consider the total momentum of a fluid element  $V$ ,  $V \subset N$ ,  $\mathbf{P}$ :

$$(177) \quad \mathbf{P} = \int_V \mathbf{p} \text{vol}^d = \int_V p^i \text{vol}^d \otimes \mathbf{e}_i$$

where  $\text{vol}^d$  is the volume  $d$ -form on the *spatial* (smooth) manifold  $N$ ,  $N$  is of dimension  $\dim N = d$ ,  $\mathbf{e}_i$  is an orthonormal basis vector of the orthonormal frame on  $N$ , and where  $i = 1, 2, \dots, d$ , with repeated indices imply summation (i.e. Einstein’s summation notation).  $\mathbf{p} = \mathbf{p}(t, \mathbf{x})$  is a time-dependent vector field taking values on the tangent bundle of  $N$ ,  $TN$ . Note that  $\mathbf{P}$  is also a time-dependent vector field, after integration of  $\mathbf{p} \text{vol}^d$  over  $V$ : strictly speaking, one could say that  $\mathbf{p}$  and  $\mathbf{P}$  belongs to the space of all vector fields over  $\mathbb{R} \times N$ , with  $\mathbf{R}$  representing time,  $\mathfrak{X}(\mathbb{R} \times N) = T(\mathbb{R} \times N)$ , but with the component in the time direction being strictly zero.

Consider the (total) time derivative of  $\mathbf{P}$ :

$$\begin{aligned} (178) \quad \dot{\mathbf{P}} &\equiv \frac{d}{dt} \mathbf{P} = \frac{d}{dt} \int_V p^i \text{vol}^d \otimes \mathbf{e}_i = \int_V \left( \frac{\partial p^i}{\partial t} + \text{div}(p^i \mathbf{u}) \right) \text{vol}^d \otimes \mathbf{e}_i = \\ &= \int_V T^{ij} dS_j \otimes \mathbf{e}_i = - \int_V (\text{grad} p) \text{vol}^d \end{aligned}$$

where the usual “right-hand side” (RHS) for the Euler equations for inviscid flow is the second line of Eq. 178, and there,  $p$  is the pressure for the thermodynamic system at  $V$ . This thermodynamic system is the (single) gas phase inside  $V$ . Take  $V$  small enough such that the system equilibrates *locally* and thus the pressure  $p$  is a good parameter to specify the thermodynamic state of the system at local equilibrium in  $V$ .  $p$  takes on double duty as a  $C^\infty$  function on  $\mathbb{R} \times N$ , i.e.  $p \in C^\infty(\mathbb{R} \times N)$ , with  $C^\infty(\mathbb{R} \times N)$  denoting the space of all smooth (i.e. infinitely differentiable) functions over (i.e. depending upon)  $\mathbb{R} \times N$ .

One thing to note about the derivation of the usual “left-hand side” (LHS) for the Navier-Stokes equations, as shown in the first line of Eq. 178, is that the Lie derivative is only acting on the “differential form” part of the integrand  $\mathbf{p} \text{vol}^d = p^i \text{vol}^d \otimes \mathbf{e}_i$ . This could be made clearer in considering the definition of a time derivative as a limit, and the flow of a differential form under the “full” velocity vector field  $\frac{d}{dt} = \frac{\partial}{\partial t} + \mathbf{u}$ ,  $\varphi_{\Delta t}$ , where  $\mathbf{u} \in \mathfrak{X}(\mathbb{R} \times N) = T(\mathbb{R} \times N)$  is the usual time-dependent velocity vector field, taking only values in  $TN$ . Indeed

$$\begin{aligned} (179) \quad \frac{d}{dt} \int_V p^i \text{vol}^d \otimes \mathbf{e}_i &:= \lim_{\Delta t \rightarrow 0} \frac{\int_{\varphi_{\Delta t} V} p^i(t + \Delta t, \mathbf{x}) \text{vol}^d \otimes \mathbf{e}_i - \int_V p^i(t, \mathbf{x}) \text{vol}^d \otimes \mathbf{e}_i}{\Delta t} = \\ &= \lim_{\Delta t \rightarrow 0} \left\{ \int_V \varphi_{\Delta t}^* (p^i(t + \Delta t, \mathbf{x}) \text{vol}^d) \otimes \mathbf{e}_i - \int_V p^i(t, \mathbf{x}) \text{vol}^d \otimes \mathbf{e}_i \right\} \frac{1}{\Delta t} = \\ &= \int_V \lim_{\Delta t \rightarrow 0} [\varphi_{\Delta t}^* (p^i(t + \Delta t, \mathbf{x}) - p^i(t, \mathbf{x}))] \frac{1}{\Delta t} \otimes \mathbf{e}_i = \\ &= \int_V \mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}} (p^i(t, \mathbf{x}) \text{vol}^d) \otimes \mathbf{e}_i \end{aligned}$$

Note that the derivation in Eq. 179 resolves the (possible) pitfall of perhaps taking the Lie derivative  $\mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}}$  over the *entire* integrand,  $\mathbf{p} \text{vol}^d \in \Omega^d(N, TN)$ , with  $\Omega^d(N, TN)$  denoting the space of all  $TN$ -valued  $d$ -forms. If that was the case, taking  $\mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}}$  over the entire integrand, we’d expect the Lie derivative to act as a “symmetric” derivation over a symmetric tensor product, i.e.  $\mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}}(\mathbf{p} \text{vol}^d) = (\mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}}(\mathbf{p})) \text{vol}^d + \mathbf{p} \otimes (\mathcal{L}_{\frac{\partial}{\partial t} + \mathbf{u}} \text{vol}^d)$ . However, the derivation in Eq. 179 shows that this is not the case.

The point is, or the take away is, that the total time derivative  $\frac{d}{dt}$  can be “moved into” the integral of a vector-valued differential form and “inside the integral,” the Lie derivative is performed.

We should take care of how mass density is defined in a small volume  $V$ . Clearly, to obtain the total mass  $M$  inside a volume  $V$ , we integrate over the mass density  $\rho$ , over  $V$ :

$$(180) \quad M = \int_V \rho \text{vol}^d$$

Clearly, when  $\rho = 0$ , then  $M = 0$ , with the physical interpretation, clearly, being that there’s no gas inside  $V$ .

Taking the total time derivative of  $M$ :

$$(181) \quad \begin{aligned} \dot{M} &= \frac{d}{dt} \int_V \rho \text{vol}^d = \int_V \frac{\partial \rho}{\partial t} \text{vol}^d + \int_V \mathbf{d}i_{\mathbf{u}} \text{vol}^d = \int_V \frac{\partial \rho}{\partial t} \text{vol}^d + \int_{\partial V} i_{\mathbf{u}} \rho \text{vol}^d = \\ &= \int_V \left( \frac{\partial \rho}{\partial t} + \text{div}(\rho \mathbf{u}) \right) \text{vol}^d \end{aligned}$$

In the first line of Eq. [181](#), the last term,  $\int_{\partial V} i_{\mathbf{u}} \rho \text{vol}^d$  expresses the mass flux through the surface of volume  $V$ ,  $\partial V$ , and then, applying Stoke’s law, leads to the second line of Eq. [181](#), leading to the divergence of  $\rho \mathbf{u}$ . Multiplying vector fields by  $C^\infty$  functions, in this case  $\mathbf{u}$  and  $\rho$ , respectively, is ok since the space of vector fields  $\mathfrak{X}(\mathbb{R} \times N)$  is a module over the commutative algebra of  $C^\infty(\mathbb{R} \times N)$ .

Mass conservation dictates that

$$(182) \quad \dot{M} = 0$$

This statement can be adapted to multiple number of species.

If  $V$  can be made arbitrary and sufficiently small, then the usual statement, in differential form, of mass conservation is obtained:

$$(183) \quad \frac{\partial \rho}{\partial t} + \text{div}(\rho \mathbf{u}) = 0 \text{ or } \frac{\partial \rho}{\partial t} = -\text{div}(\rho \mathbf{u})$$

Now, one should take special care in defining the “momentum density”, the momentum per unit volume,  $\mathbf{p}$ , with respect to the mass density, given a velocity vector field  $\mathbf{u}$ :

$$(184) \quad \mathbf{p} = \rho \mathbf{u} \text{ or } p^i = \rho u^i \quad i = 1, \dots, d$$

If  $\rho = 0$ , then  $\mathbf{p} = 0$ . But it can be that  $\mathbf{u} \neq 0$ , say, at a particular point in  $(t, \mathbf{x})$ . Then  $\frac{\mathbf{p}}{\rho}$  is ill-defined.

If  $\mathbf{u} = 0$ , but  $\rho > 0$ , then the physical interpretation is clear: the gas is at rest.  $\mathbf{p} = \mathbf{0}$  is thus clear. Then  $\frac{\mathbf{p}}{\rho}$  is well-defined and so is  $\mathbf{u}$  at that point  $(t, \mathbf{x})$ .

So Eq. [183](#) gives a direct recipe for the use of finite difference methods in the next time step  $t + \Delta t$ , where  $\Delta t$  is infinitesimally small, for  $\rho$ ,  $\rho(t + \Delta, \mathbf{x})$ . But the important point is that we may not be able to define what  $\mathbf{u}$  will be in the next time step.

From Eqns. [178](#), [179](#), we obtain a direct recipe for the use of finite difference methods in the next time step  $t + \Delta t$  for  $\mathbf{p}$  or  $p^i$ ,  $i = 1 \dots d$ :

$$(185) \quad \frac{\partial p^i}{\partial t} + \text{div}(p^i \mathbf{u}) = -(\text{grad} p)^i \text{ or } \frac{\partial p^i}{\partial t} = -\text{div}(p^i \mathbf{u}) - (\text{grad} p)^i$$

As a notation note,  $p$ , that isn’t decorated by superscript  $i$ ’s or isn’t boldfaced, refers to the *pressure*, indeed the “local” pressure at point  $(t, \mathbf{x})$ ;  $p^i$  or  $\mathbf{p}$  refer to the momentum density.

So from Eq. [185](#), we obtain  $\mathbf{p}(t + \Delta t, \mathbf{x})$ , but we don’t obtain  $\mathbf{u}$  directly. If  $\rho > 0$  (and the fact that  $\rho$  is non-negative is a *physical* fact that must be enforced in the math), then dividing  $\mathbf{p}$  by  $\rho$  yields  $\mathbf{u}$ :

$$(186) \quad \frac{\mathbf{p}}{\rho} = \mathbf{u} \quad \text{if } \rho > 0$$

However, if  $\rho = 0$ , then  $\mathbf{u}$  is undetermined.

My concerns are the following: in the implementation of finite difference methods, what is “small enough” for  $\rho$  to warrant it to be effectively zero? In Eq. [186](#), suppose  $\rho$  is very small, but due to the floating point representation of numbers by the computer away from 0, then  $\mathbf{u}$  could be very large, perhaps “overflowing” the floating pointer presentation of the number.

On one level, I would think that we can check if  $\rho$  is equal to 0 or not and then if not, calculate the new  $\mathbf{u}(t + \Delta t, \mathbf{x})$ . If  $\rho$  is equal to 0, then  $\mathbf{u}$  is left alone.

**27.6. Energy density.** Consider a thermodynamic system at rest. Indeed, consider moving into the reference frame where the fluid element is at rest inside the volume  $V$ ; this can be done at any instant  $t$ .

The enthalpy  $H$  of this system is defined as

$$(187) \quad H := U + pV$$

with  $U$  denoting the *internal energy* of the thermodynamic system.  $U$  includes all the energy locked up in rotational and vibrational modes, chemical potential, random thermal kinetic energy, but *not* the kinetic energy due to the *bulk* fluid motion due to velocity vector field  $\mathbf{u}$ .

$H$  is a function of entropy (entropy that’s *not* multiplied by the Boltzmann constant; it is dimensionless in this notation)  $\sigma$  of the system, pressure  $p$ , and the number of molecules of the species; and if there’s only one species present, this is denoted as  $N$ ; otherwise, for each species  $i$ ,  $H$  is a function of  $\{N_i\}_i$ , the number of each species present:

$$H = H(\sigma, p, N) \text{ or } H = H(\sigma, p, \{N_i\}_i)$$

In fact, the parameters  $\sigma, p, \{N_i\}_i$  or  $(\sigma, p, N)$  in the single-species case, define exactly (i.e. that’s all you need to know) the state of the thermodynamic system at a given time and place  $(t, \mathbf{x})$  - let these parameters be global coordinates for a manifold  $\Sigma$  that represents all possible thermodynamic states of the system.

$U$  is a function of  $\sigma, V, N$  or  $\sigma, V, \{N_i\}_i$  in the multi-species case, which also specify completely the thermodynamic state of the system

$$U = U(\sigma, V, N)$$

Going from  $U$  to  $H$  is a Legendre transformation on  $\Sigma$  -  $\Sigma$  stays the same, but the (global) coordinates are transformed.

Nevertheless, when there is no gas in  $V$ , i.e.  $N = 0$ , then define the zero  $U = 0$  for  $N = 0$ .

$$U = U(\sigma, V, 0) = U(0, V, 0) = 0$$

Notice that  $\sigma = 0$  in the above expression. The physical argument is this: if there is no gas in the box of volume  $V$ , there is only 1 way to arrange the microstates of the thermodynamic system in such a manner, and so the log of 1 is 0 and so  $\sigma = 0$ .

For  $H$ , we have to define what it means for  $H$  to be zero carefully. If  $N = 0$ , then  $U = 0$  and so

$$H(\sigma, p, 0) = U(\sigma, V, 0) + pV = 0 + pV$$

I will argue that  $p = 0$  if  $N = 0$ . For if  $N = 0$ , then there’s no pressure to push back against the target gas system. This is clearly the case of a vacuum. For the case where there’s already an ambient or background gas phase in and around  $V$ , this is also the case as well as we’re concerned with the target gas system. The background gas can be pushing against itself and have a pressure inside and around  $V$ , or even a pressure gradient, but if there’s no target gas there at  $V$ , then there’s nothing for this background gas to push against, other than itself.

So  $H = 0$  when  $N = 0$ .

Define the (internal) energy density in this reference frame (denoted with  $'$ ) where the fluid is at rest, the (internal) energy per unit volume,  $\epsilon'$  to be

$$(188) \quad \epsilon' := \frac{U}{V}$$

Define the enthalpy per unit volume  $h'$  to be

$$(189) \quad h' := \frac{H}{V} = \epsilon' + p$$

If  $\epsilon' = 0$ , and yet temperature of the thermodynamic system  $\tau$  is positive, i.e.  $\tau > 0$ , then no gas is there. Thus, also  $h' = 0$  as  $p = 0$  as argued above.

As a notation note,  $\tau$  denotes the temperature in units of *energy* such as joules ( $J$ ) or (ergs): use Boltzmann’s constant  $k_B$  as a units conversion factor between  $\tau$  and the usual temperature  $T$  in Kelvin

$$\tau = k_b T$$

Why I’ve elaborated upon when  $U = 0$  and  $H = 0$  and thus when  $\epsilon' = 0$  and  $h' = 0$  is because, first, the energy scale “zero” can be arbitrarily set and so it’s good to set it once and for all, with physically meaningful reasons. Second, there can be confusion that arises in the consideration of what it means for  $H = 0$  as when it is defined as  $H := U + pV$ , and of what pressure  $p$  refers to: the pressure that the target gas system encounters, only, or would it include the “ambient” pressure even if there’s no gas in the local volume  $V$  we’re considering. Third, these considerations stem from the concern of what it means to, in a volume  $V$  with no gas initially present, put 1 molecule of gas into the fixed volume  $V$  and the resulting changes in thermodynamic quantities; for instance, where there was no gas before, putting in 1 gas molecule results immediately in  $pV$

work that needs to be done to restore the volume back to its initial “fixed” volume  $V$  (see Le Bellac, Mortessagne, Batrouni (2004) [15] and its lucid distinguishing between *convection* and *conduction* for a thermodynamic system).

Clearly, when switching, i.e. “boosting” into the reference frame where the fluid is moving, i.e. the observer’s or laboratory’s frame of reference,  $\epsilon'$  transforms into  $\epsilon$ , the (internal) energy density in this lab frame, in the following manner:

$$(190) \quad \epsilon = \epsilon' + \frac{1}{2}\rho u^2$$

Denote this kinetic energy density  $\frac{1}{2}\rho u^2$  as  $k \equiv \frac{1}{2}\rho u^2$ .

For total energy in a volume  $V$ ,  $E$ , given by

$$(191) \quad E := \int_V \epsilon \text{vol}^d$$

then the total time derivative of  $E$ ,  $\dot{E} \equiv \frac{dE}{dt}$  is given by, from differential geometry as above,

$$(192) \quad \begin{aligned} \dot{E} &\equiv \frac{d}{dt} \int_V \epsilon \text{vol}^d = \int \left( \frac{\partial \epsilon}{\partial t} + \text{div}(\epsilon \mathbf{u}) \right) \text{vol}^d = \\ &= \int \left( \frac{\partial \epsilon}{\partial t} + \text{div}(k + \epsilon') \mathbf{u} \right) \text{vol}^d = \int \left( \frac{\partial \epsilon}{\partial t} + \text{div}(k \mathbf{u}) + \text{div}((h' - p) \mathbf{u}) \right) \text{vol}^d = \int_V \left( \frac{\partial \epsilon}{\partial t} + \text{div}(k \mathbf{u}) + \text{div}(h' \mathbf{u}) - \text{div}(p \mathbf{u}) \right) \text{vol}^d \end{aligned}$$

From physics, the total time derivative of the energy  $E$ ,  $\dot{E}$  is equal to the work done on the fluid system inside and about  $V$  by external forces. Neglecting body forces, the only external forces are due to stresses and strains at the surface of  $V$ , which are wrapped up in the stress-strain tensor  $T^{ij}$ . Using Stoke’s equation again, a surface integral involving  $T^{ij}$  can be transformed into a volume integral. For the case of the perfect fluid,  $T^{ij}$  takes the form of  $T^{ij} = g^{ij}(-p)$ , where  $g^{ij}$  is the metric of the spatial manifold  $N$ . Thus

$$(193) \quad \dot{E} = \int_V \text{div}(u_i T^{ij}) \text{vol}^d = \int_V \text{div}(u_i g^{ij}(-p)) \text{vol}^d = - \int_V \text{div}(\mathbf{u} p) \text{vol}^d$$

Recognizing that we can subtract from both sides the  $-\int_V \text{div}(\mathbf{u} p) \text{vol}^d$  term of the “left-hand side” and “right-hand side” of the “energy balance” equation for  $\dot{E}$ , Eqns. [192](#), [193](#), respectively, then we obtain the desired integral and differential form for the energy density  $\epsilon$  under compressible fluid flow:

$$(194) \quad \implies \int_V \left( \frac{\partial \epsilon}{\partial t} + \text{div}((k + h') \mathbf{u}) \right) \text{vol}^d = 0 \implies \boxed{\frac{\partial \epsilon}{\partial t} + \text{div}((k + h') \mathbf{u}) = 0}$$

Now if we can express  $h'$  in terms of  $\rho$  and  $\mathbf{u}$ , given the concerns above, and  $\epsilon$ , then we’d have succeeded in a complete or “closed” set of differential equations for dynamically evolving the fluid system.

Consider using the heat constants for constant volume and constant pressure,  $C_V$ ,  $C_p$ , respectively. Definitively or exactly,

$$dU = C_V N d\tau$$

$$dH = C_p N d\tau$$

These statements are exact.

Assuming that  $C_V$  and  $C_p$  are constant over the range of temperatures we’re interested in (note that this case can be modified as  $C_V$ ,  $C_p$  can be interpolated from experimentally obtained values and so they’re *not* constant), then

$$dU = C_V N d\tau \xrightarrow{f} U = C_V N \tau$$

where for  $\tau = 0$ ,  $U = 0$  (I’ve set the zero for  $U$  to be such). Also notice that when  $N = 0$ ,  $U = 0$  - physically, we’d expect that if there’s no gas there at  $V$ , then  $U = 0$ .

Also

$$\int_{\gamma} dH = \int_{\gamma} C_p N d\tau = C_p N (\tau - \tau_0) = H - H_0$$

where  $\gamma$  is a curve on  $\Sigma$ .

Let’s set the zero for  $H$ . For  $\tau_0 = 0$ , consider

$$H - H_0 = C_p N \tau \implies H = C_p N \tau$$

Now

$$H_0 = U_0 + pV = 0$$

where  $pV = N\tau$ , the ideal gas law, was used for  $\tau = 0$ , as for  $\tau = 0$ ,  $p = 0$  if  $V \neq 0$ .

Now that the zeros for  $H$  and  $U$  are set, then we can definitely say that

$$H = U + pV = C_V N \tau + pV = C_p N \tau \text{ or } U = C_p N \tau - pV$$

Now

$$\gamma := \frac{C_p}{C_V}$$

and so

$$\gamma - 1 = \frac{C_p}{C_V} - 1 = \frac{C_p - C_V}{C_V}$$

From  $H = U + pV$ , and the expressions above,

$$C_p N \tau = C_V N \tau + pV \text{ or } (C_p - C_V) N \tau = pV$$

$$\implies (\gamma - 1) C_V N \tau = (\gamma - 1) U = pV$$

$$\xrightarrow{\frac{1}{V}} (\gamma - 1) \epsilon' = p$$

$$\implies \boxed{p = (\gamma - 1) \left( \epsilon - \frac{1}{2} \rho u^2 \right)} \text{ or } p = (\gamma - 1) (\epsilon - k)$$

Thus, we can reexpress  $h'$  in a number of ways using the expression above:

$$h' = \epsilon' + p = \epsilon' + (\gamma - 1) \epsilon' = \gamma \epsilon'$$

Thus,

$$\frac{\partial \epsilon}{\partial t} + \text{div}(k \mathbf{u}) + \text{div}(\gamma \epsilon' \mathbf{u}) = \frac{\partial \epsilon}{\partial t} + \text{div}(k \mathbf{u}) + \text{div}(\gamma (\epsilon - k) \mathbf{u}) =$$

$$\frac{\partial \epsilon}{\partial t} + \text{div}(k \mathbf{u} (1 - \gamma)) + \text{div}(\gamma \epsilon \mathbf{u}) = 0$$

where  $\epsilon' = \epsilon - k$  was used.

Thus

$$(195) \quad \boxed{\frac{\partial \epsilon}{\partial t} = -\text{div}(\gamma \epsilon \mathbf{u}) - \text{div}(k \mathbf{u} (1 - \gamma))} \text{ where } k := \frac{1}{2} \rho u^2$$

I want to point out that an advantage of this expression for  $\epsilon$ , Eqn. [195](#), is that there is no division so we are not worried about dividing by 0; this is a good trait to have in numerical computation, and on the theoretical side, division by  $C^\infty$  functions is undefined for the vector field as a section of the tangent bundle  $TN$ .

Of note, it should be remarked that the temperature  $\tau$  can be easily obtained once  $\epsilon$  and  $\mathbf{u}$  is known at each point  $(t, \mathbf{x})$ :

$$\epsilon = k + \epsilon' = k + c_V \tau$$

$$\implies \boxed{\tau = \frac{\epsilon - k}{c_V}}$$

where  $c_V := \frac{C_V}{V}$  the heat capacity at constant volume, per unit volume.

To obtain  $\epsilon$  over  $N$ , as is necessary for our numerical computation to work, recall the following definitions and following steps:

$$(196) \quad \begin{aligned} \epsilon &:= \epsilon' + \frac{1}{2}\rho u^2 & \rho \geq 0 \\ \epsilon' &:= \frac{U}{V} = \frac{C_V N \tau}{V} = \frac{C_V \rho \tau}{M} \\ &\boxed{\epsilon = \rho \left( \frac{C_V \tau}{M} + \frac{1}{2}u^2 \right)} \end{aligned}$$

Again, I want to emphasize the advantage of the last expression in Eq. 196 in that it’s clear, in the physical interpretation, that when  $\rho = 0$ , no gas is present at that point  $(t, \mathbf{x})$ , then  $\epsilon = 0$  and that there are no division arithmetic, avoiding division error in numerical computation and the fact that the division operator is undefined on the theory side (that the space of vector fields,  $\mathfrak{X}(\mathbb{R} \times N)$  is a R-module).

Thus, here is the strategy for the numerical computation for compressible fluid flow for the case of the ideal gas:

Given  $\rho$ , initially, everywhere on  $N$ ,

Given an initial  $\mathbf{u}$  everywhere on  $N$ ,

$\mathbf{p}$  is easily obtained through  $\mathbf{p} = \rho \mathbf{u}$ .

Given an initial temperature “distribution”  $\tau$ , then we have an initial “distribution” for what  $\epsilon$  is at each point in  $N$ , through

$$\epsilon = \rho \left( \frac{C_V \tau}{M} + \frac{1}{2}u^2 \right)$$

we can then use the following complete set of differential equations that dynamically evolves the system at each small volume  $V$  is as follows:

$$(197) \quad \boxed{\begin{aligned} \frac{\partial \rho}{\partial t} &= -\text{div}(\rho \mathbf{u}) \\ \frac{\partial p^i}{\partial t} &= -\text{div}(p^i \mathbf{u}) - (\text{grad} p)^i \\ \frac{\partial \epsilon}{\partial t} &= -\text{div}(\gamma \epsilon \mathbf{u}) - \text{div}(k \mathbf{u}(1 - \gamma)) \end{aligned}}$$

with the pressure  $p$  given by

$$p = (\gamma - 1)\left(\epsilon - \frac{1}{2}\rho u^2\right)$$

at all points  $(t, \mathbf{x})$ .

The only adhoc assertion is obtaining the velocity vector field at the new time step, denoted  $\mathbf{u}(t + \Delta t, \mathbf{x})$ :

$$\frac{\mathbf{p}(t + \Delta t, \mathbf{x})}{\rho(t + \Delta t, \mathbf{x})} = \mathbf{u}(t + \Delta t, \mathbf{x})$$

if  $\rho(t + \Delta t, \mathbf{x}) \neq 0$ .

Otherwise,  $\mathbf{u}(t + \Delta t, \mathbf{x}) = u(t, \mathbf{x})$ .

Another possibility for  $\mathbf{u}$  is if, in the absence of the target gas system, that it obeys the geodesic equation. If this is the case, then for the 4-vector that is the “4-velocity”,  $u \in \mathfrak{X}(\mathbb{R} \times N)$ , which takes on 2 different forms (and we’ll use the nonrelativistic version)

$$\begin{aligned} u &= (-1, \beta) && \text{relativistic} \\ u &= (1, \mathbf{u}) && \text{nonrelativistic form} \end{aligned}$$

$$\begin{aligned} \nabla_u u &= u^\rho \nabla_{\frac{\partial}{\partial x^\rho}} u = u^\rho \left\{ \frac{\partial u}{\partial x^\rho} + \Gamma_{\nu\rho}^\mu u^\nu \frac{\partial}{\partial x^\mu} \right\} = \\ &= \frac{\partial \mathbf{u}}{\partial t} + u^k \frac{\partial \mathbf{u}}{\partial x^k} + \Gamma_{00}^\mu \frac{\partial}{\partial x^\mu} + \Gamma_{k0}^\mu u^k \frac{\partial}{\partial x^\mu} + u^k \Gamma_{0k}^\mu \frac{\partial}{\partial x^\mu} + u^k \Gamma_{jk}^\mu u^j \frac{\partial}{\partial x^\mu} = \\ &= \frac{\partial \mathbf{u}}{\partial t} + u^k \frac{\partial \mathbf{u}}{\partial x^k} + \Gamma_{00}^\mu \frac{\partial}{\partial x^\mu} + 2\Gamma_{k0}^\mu u^k \frac{\partial}{\partial x^\mu} + u^k u^j \Gamma_{jk}^\mu \frac{\partial}{\partial x^\mu} \end{aligned}$$

## 28. HAMILTONIAN

Novikov (1982) [24] gives the Hamiltonian formalism for a compressible ideal fluid.

In Sec. 2 “The Hamiltonian Formalism of Systems of Hydrodynamic Origin” of Novikov (1982) [24], the Hamiltonian formalism “usually” (my words) involves “Clebsch variables” but Novikov points out that these field variables can’t always be introduced, and if so, then frequently only locally. Novikov advocates the language of “Poisson brackets” over the language of symplectic manifolds (of even dimensions). One reason is that the number of field variables is odd.

## 29. HEAT EQUATION

**29.1. Boundary Conditions.** Let  $N = \mathbb{R}^3$ . Let  $B \subset \mathbb{R}^3$ ,  $B := \{(x, y, z) \in \mathbb{R}^3 | 0 \leq x \leq l_x, 0 \leq y \leq l_y, 0 \leq z \leq l_z\}$ ;  $(l_x, l_y, l_z) \in (\mathbb{R}^{+^3})$ , i.e.  $\mathbb{R}^+ \subset \mathbb{R}$ , i.e. all positive real numbers (still an open set), with  $l_x, l_y, l_z > 0$ .

Consider the “inlet” hypersurface at  $x = 0$ ,  $S_{\text{inlet}}$ , so  $S_{\text{inlet}} \subset B$  (in fact,  $S_{\text{inlet}} \subset \partial B$ ,  $\dim S_{\text{inlet}} = d - 1 = 3 - 1 = 2$ . Define as such:

$$(198) \quad S_{\text{inlet}} := \{(x, y, z) \in B | x = 0\}$$

Consider the “sides” hypersurface,  $S_{\text{sides}}$ ,  $S_{\text{sides}} \subset B$ , in fact,  $S_{\text{sides}} \subset \partial B$ .

Now

$$(199) \quad S_{\text{sides}} := S_{y=0} \coprod S_{y=l_y} \coprod S_{z=0} \coprod S_{z=l_z}$$

with

$$(200) \quad \begin{aligned} S_{y=0} &= \{(x, y, z) \in B | y = 0\} \\ S_{y=l_y} &= \{(x, y, z) \in B | y = l_y\} \\ S_{z=0} &= \{(x, y, z) \in B | z = 0\} \\ S_{z=l_z} &= \{(x, y, z) \in B | z = l_z\} \end{aligned}$$

## 30. LATTICE BOLTZMANN METHOD

Pitaevskii, Lifshitz’s **Physical Kinetics** (1981) [30]

Bao and Meskas (2011) [25]

Tölke (2008) [31]

A reference for the so-called  $d2q9$  model is given by Qian, D’Humières, and Lallemand (1991) [32].

In the Lattice Boltzmann method, the velocity space (really a  $A^d$  affine space) is discretized into a discrete set of microscopic velocities  $\{\mathbf{e}_i\}_{i \in I}$ , for the  $i$ th particle.

For the model

$D2Q9 \equiv d2q9$

$D3Q15 \equiv d3q15$

$D3Q19 \equiv d3q19$

respectively (wikipedia and Qian’s notation vs. Tölke’s notation), the microscopic velocities are (according to wikipedia)

(201)

$$\mathbf{e}_i = c \times \begin{cases} (0, 0) & i = 0 \\ (1, 0), (0, 1), (-1, 0), (0, -1) & i = 1, 2, 3, 4 \\ (1, 1), (-1, 1), (-1, -1), (1, -1) & i = 5, 6, 7, 8 \end{cases}$$
$$\mathbf{e}_i = c \times \begin{cases} (0, 0, 0) & i = 0 \\ (\pm 1, 0, 0), (0, \pm 1, 0), (0, 0, \pm 1) & i = 1, 2, \dots, 5, 6 \\ (\pm 1, \pm 1, \pm 1) & i = 7, 8, \dots, 13, 14 \end{cases}$$
$$\mathbf{e}_i = c \times \begin{cases} (0, 0, 0) & i = 0 \\ (\pm 1, 0, 0), (0, \pm 1, 0), (0, 0, \pm 1) & i = 1, 2, \dots, 5, 6 \\ (\pm 1, \pm 1, 0), (\pm 1, 0, \pm 1), (0, \pm 1, \pm 1) & i = 7, 8, \dots, 17, 18 \end{cases}$$

31. REACTIVE FLOWS

31.1. **Kyle E. Niemeyer - Kyle E. Niemeyer’s stuff.** Niemeyer (2013) [26]

32. SPARSE MATRICES

32.1. **Compressed Sparse Row (CSR).** Consider  $M \times N$  matrix  $A$ . Let  $|A|$  = total number of nonzero entries in  $M$ ,  $|A| < MN$ .

Let  $V = \{A_{ij}|A_{ij} \neq 0\}$ , so  $|V| = |A|$ .

Let  $J \equiv$  column indices,  $|J| = |A|$ , consist of column index  $j_k$  of  $k$ th non-zero entries.

So

$$V_k = A_{i_k j_k} \text{ s.t. } A_{i_k j_k} \neq 0, k = 0, 1, \dots |A| - 1$$
$$J_k = j_k, \quad k = 0, \dots, |A| - 1$$

Let  $I \equiv$  row indices,  $|I| = M + 1$ .  
It encodes the index in  $V, J$ , where the given row starts; i.e.  
 $I[j]$  encodes total number of nonzero elements above row  $j$ .  
 $I[M] \equiv I_M = |A| \equiv$  total number of nonzero entries.

So  $I_l \in \{0, 1, \dots |A|\}$ .

Note that  $I_{l+1} - I_l =$  total number of nonzero elements in row  $l$ .

Consider  $k$ th nonzero elements,  $k = 0, 1, \dots |A| - 1$ .  
 $V_k$  is its value.  
 $J_k = j_k$

For  $I_l, l = 0, 1, \dots M$ , find  $l$  s.t.  $k \geq I_l$  and  $I_{l+1} < k$ .

Thus  $A_{ij} = V_k$  with  $j = j_k$  and  $i = l$  s.t.  $k \geq I_l$  and  $I_{l+1} < k$ .  
In summary,

**Definition 9** (Compressed Sparse Row (CSR)). *For  $M \times N$  matrix  $A$ ,  $|A|$  = total number of nonzero entries in  $M$ , Let  $V = \{A_{ij}|A_{ij} \neq 0\}$ ,  $|V| = |A|$ ,  $J \equiv$  column indices,  $|J| = |A|$ ,*

*$I \equiv$  row offsets,  $|I| = M + 1$ ,  
 $I[l] =$  total number of nonzero elements in rows  $0, 1, \dots l - 1$ ,  $I[M] = |V|$*

*For  $k = 0, 1, \dots |A| - 1$ ,  
 $A_{i_k j_k} = V_k$  ,  $j_k = J[k]$  and  $i_k = l$  s.t.  $I[l + 1] > k$  and  $I[l] \leq k$*

*Given  $i, j$ ,  $I[i] \leq k < I[i + 1]$   
Then for  $k = I[i], \dots I[i + 1] - 1$ , try  $J[k] = j$ . Once  $k$  found  $A_{ij} = V_k$ .*

32.1.1. *Matrix multiplication with Compressed Sparse Row matrices.* Consider  $A_{ij}x_j = y_i$ . Consider  $k = I[i], \dots I[i + 1] - 1$ , i.e. consider  $I[i + 1] - I[i]$  nonzero elements in  $i$ th row.

$$A_{ij}x_j = \sum_{k=I[i]}^{I[i+1]-1} V_kx[J[k]]$$

Since  $J[k] = j_k =$  column index of the  $k$ th element.

32.1.2. *Example: Tridiagonal symmetric matrix.* Consider  $M \times N$  matrix  $A$ , s.t.  $i = 0, 1, \dots M - 1, j = 0, 1, \dots N - 1$ .

Diagonal entries  $A_{ii}$ .  $i = 0, 1, \dots \min \{M, N\}$ .

Consider  $A_{ij}$  s.t.  $j = i + 1$ . If  $M = N$ ,  $i = 0, 1, \dots M - 2$ , and so  $M - 1$  entries.

If  $M > N$ ,  $j = 1, \dots N - 1$ , so  $N - 1$  entries.  
If  $M < N$ ,  $i = 0, \dots M - 1$ , and so  $\min \{M - 1, N - 1\}$  entries.

Consider  $A_{ij}$  s.t.  $i = j + 1$ . Likewise  $\min \{M - 1, N - 1\}$  entries.

Part 8. Optimization

33. MINIMIZATION OR MAXIMIZATION OF FUNCTIONS

Ch. 10. "Minimization or Maximization of Functions", Press, Teukolsky, Vetterling, Flannery (2007) [36]

33.1. **Initially Bracketing a Minimum.** Sec. 10.1 "Initially Bracketing a Minimum", Press, Teukolsky, Vetterling, Flannery (2007) [36]

33.1.1. *Pseudocode.* Pseudocode for "Bracket", for initial bracketing of a minimum. The purpose is to show the existence of a minimum.  
 $G_l = 1.618032$ ,  
 $G_L = 100.0$   
 $T = 1 \times 10^{-20}$

Let  $a, b \in \mathbb{R}$ ,  
Consider  $f : \mathbb{R} \rightarrow \mathbb{R}$ , and  $f(a) \geq f(b)$ . Otherwise for  $f(a) < f(b)$ , then switch  $a, b$ .

$$c = b + G_l(b - a)$$



While  $f(b) > f(c)$ ,

$$\begin{aligned} r &= (b-a)(f(b)-f(c)) \\ q &= (b-c)(f(b)-f(a)) = (c-b)(f(a)-f(b)) \\ u &= b - ((b-c)*q - (b-a)r)/(2\max\{|q-r|, T\} \frac{q-r}{|q-r|}) \\ u_{\text{lim}} &= b + G_L(c-b) \end{aligned}$$

If  $(b-u)(u-c) > 0$ , i.e.  $b < u < c$  or  $b > u > c$  (and  $f(a) > f(b)$ ),

$$\begin{aligned} fu &= f(u) \\ \text{if } fu < f(c), & \text{ go a minimum between } b \text{ and } c \\ b, u, c \end{aligned}$$

else if  $f(b) < f(u)$ , got a minimum between  $a$ , and  $u$   
return  $a, b, u$

Otherwise,  $u$  is between  $b$  and  $c$  but  $f(b) \geq f(u) \geq f(c)$ , and so  
parabolic fit was no use. Use default magnification:  
 $u = c + G_l(c-b)$   
 $fu = f(u)$

else if  $(c-u)(u-u_l) > 0.0$  i.e.  $c < u < u_l$  or  $(c > u > u_l)$  (and also  $(b-u)(u-c) \leq 0.0$ )  
i.e. Parabolic fit is between  $c$  and its allowed limit. if  $f(u) < f(c)$   
shift  $b_x, c_x, u, u + G_l(u - c_x)$   
shift  $f(b), f(c), f(u), f(u)$

else if  $(u_l - u)(c - u_l) > 0.0$ , i.e. limit parabolic  $u$  to maximum allowed value,  
 $u := u_l$   
 $f_u := f(u)$

else  $(u_l - u)(c - u_l) \leq 0.0$ , and  $(c - u)(u - u_l) \leq 0$ ,  
reject parabolic  $u$ , use default magnification.

$$\begin{aligned} u &= c + G_l(c-b) \\ fu &= f(u) \end{aligned}$$

Otherwise,  
shift  $a, b, c, u$ ,  
shift  $f(a), f(b), f(c), f(u)$

33.2. Parabolic Interpolation, Parabolic Extrapolation. See <https://drlvk.github.io/nm/section-parabolic-interpolation.html>

Consider Newton interpolating polynomial:

$$\begin{aligned} g(x) &= f(a) + \alpha(x-a) + \beta(x-a)(x-b) \\ g'(x) &= \alpha + \beta(2x-a-b) \\ g''(x) &= 2\beta \end{aligned} \tag{202}$$

If  $g'(x) = 0$ , then  $x = \frac{a+b}{2} - \frac{\alpha}{2\beta}$ .

$$g(b) = f(a) + \alpha(b-a) \text{ so let } g(b) = f(b), \text{ so } \frac{f(b)-f(a)}{b-a} = \alpha$$

$$g(c) = f(a) + \alpha(c-a) + \beta(c-a)(c-b), \text{ so let } g(c) = f(c), \text{ so } \beta = \frac{f(c)-f(a)-\alpha(c-a)}{(c-a)(c-b)}$$

and so in summary, consider the interpolating parabola  $g(x)$  for 3 points  $a, b, c$

$$\begin{aligned} g(x) &= f(a) + \alpha(x-a) + \beta(x-a)(x-b) \text{ where} \\ \alpha &:= \frac{f(b)-f(a)}{b-a} \\ \beta &:= \frac{f(c)-f(a)-\alpha(c-a)}{(c-a)(c-b)} \text{ and} \\ x &= \frac{a+b}{2} - \frac{\alpha}{2\beta} \text{ for } g'(x) = 0 \\ g''(x) &= 2\beta \end{aligned} \tag{203}$$

Change notation:  $a \leftrightarrow b$   
 $\beta \mapsto \beta = \frac{f(c)-f(b)-\alpha(c-b)}{(c-a)(c-b)} = \frac{f(c)-f(b)}{(c-a)(c-b)} - \frac{\alpha}{c-a}$  or  
 $\beta(c-a)(c-b)(b-a) = -r - (f(b)-f(a))(c-b) = -r + q$  where

$$\begin{aligned} r &:= (b-a)(f(b)-f(c)) \\ r &= q - \beta(c-a)(c-b)(b-a) \\ q &:= (f(b)-f(a))(b-c) = (b-a)(b-c)\alpha \end{aligned}$$

If we take a look at the code for `mins.h` for `void bracket(..)` for `struct Bracketmethod` of pp. 491, Press, Teukolsky, Vetterling, Flannery (2007) [36], then  $u$  is found to be where our interpolating parabola  $g(x)$  has a zero value for the derivative, i.e.  $g'(x) = 0$ :

$$\begin{aligned} u &= b - \frac{(b-c)q - (b-a)r}{2(q-r)} = b - \frac{(b-c)q - (b-a)(q - \beta(c-a)(c-b)(b-a))}{2(q-r)} = \\ &= b - \frac{(a-c)q + \beta(c-a)(c-b)(b-a)^2}{2\beta(c-a)(c-b)(b-a)} = b - \frac{b-a}{2} + \frac{-\alpha}{2\beta} = \frac{b+a}{2} - \frac{\alpha}{2\beta} \end{aligned}$$

33.2.1. *Derivation of Initially Bracketing a Minimum code directly from Parabolic Interpolation.* Consider instead using Eq. 203 directly:

$$\begin{aligned} u &= \frac{a+b}{2} - \frac{\alpha}{2\beta} = \frac{1}{2} \left[ a+b - \frac{f(b)-f(a)}{b-a} \frac{(c-a)(c-b)}{f(c)-f(a) - \frac{f(b)-f(a)}{b-a}(c-a)} \right] = \\ &= \frac{1}{2} \left[ a+b - \frac{(f(b)-f(a))(c-a)(c-b)}{(b-a)(f(c)-f(a)) - (f(b)-f(a))(c-a)} \right] \end{aligned} \tag{204}$$

Redefine  $q, r$  to be this:

$$q := (b-a)(f(c)-f(a)) \tag{205}$$

$$r := (c-a)(f(b)-f(a)) \tag{206}$$



So then

$$(207) \quad u = \frac{1}{2} \left[ (a+b) - \frac{r(c-b)}{q-r} \right]$$

$$(208) \quad \alpha = \frac{r}{(c-a)(b-a)}$$

$$(209) \quad \beta = \frac{\frac{q}{b-a} - \frac{r}{(c-a)(b-a)}(c-a)}{(c-a)(c-b)} = \frac{q-r}{(b-a)(c-a)(c-b)}$$

And let’s rederive and reformulate the initially bracketing a minimum code.

Let  $a, b$  s.t.  $f(a) > f(b)$ . Let  $G_l > 0$ . Let  $c = b + G_l(b - a)$ .

If  $f(b) < f(c)$ , done.

Otherwise, given that  $f(b) \geq f(c)$ ,

find  $u$  s.t.  $g'(u) = 0$  and  $c < u < u_{\text{lim}}$  or  $c > u > u_{\text{lim}}$  for some limit on  $u$ ,  $u_{\text{lim}}$ .

If  $c < u_{\text{lim}} < u$  or  $c > u_{\text{lim}} > u$  then we fail finding a minimum bracket.

Otherwise,

if  $a < u < b$  or  $a > u > b$  and  $f(u) < f(b)$ , return  $(a, u, b)$ .

if  $b < u < c$  or  $b > u > c$  and  $f(b) < f(a)$  or  $f(u) < f(c)$ , return  $(b, u, c)$

if  $b < c < u$  or  $u < c < b$  and  $f(u) > f(c)$ , return  $(b, c, u)$ .

Otherwise fail.

If we had failure and  $c < u_{\text{lim}} < u$ , or  $c > u_{\text{lim}} > u$ , then let  $u = u_{\text{lim}}$ .

Otherwise, if we had failure, then let  $u = c + G_l(c - b)$ , i.e. we continue to advance in the same direction.

**33.3. Golden Section Search in 1-dim.** See Sec. 10.2 ”Golden Section Search in One Dimension”, Press, Teukolsky, Vetterling, Flannery (2007) [36]

**33.3.1. How small of a bracket about a minimum  $(a, b, c)$  is ”tolerably” small?** How small of a bracket about a minimum  $(a, b, c)$  is ”tolerably” small?

Consider Taylor’s Thm., with  $f'(b) = 0$ .

$$f(x) \approx f(b) + \frac{1}{2}f''(b)(x-b)^2$$

Want  $\epsilon|f(b)| > \frac{1}{2}|f''(b)|(x-b)^2$  or  $|x-b| < \sqrt{\epsilon}|b|\sqrt{\frac{2|f(b)|}{b^2f''(b)}}$ .  $f''(b) > 0$  for a minimum.

Included  $\frac{1}{b^2}$  because  $\sqrt{\frac{2|f(b)|}{f''(b)}} \sim 1$  (number of order unity) while  $\sqrt{\epsilon}$  is  $10^{-4}$  or  $10^{-8}$ .

Choose tolerance `tol` to be square root of machine’s floating-pt. precision.

**33.3.2. Step Derivation of Golden Section Search.** Consider problem of choosing new pt.  $x$ , given  $(a, b, c)$ .

Suppose  $\frac{b-a}{c-a} = w$ ,  $\frac{c-b}{c-a} = 1 - w$ .

Suppose additional fraction  $z$  beyond  $b$ ,  $\frac{x-b}{c-a} = z$ .

Next, bracketing segment will either be  $w + z$  or  $1 - w$  (depending on where minimum is).

Choose  $z$  to make these lengths equal, so to ”hedge our bets” equally (no prior information).

$w + z = 1 - w$  or  $z = 1 - 2w$ .

Note that  $z > 0$  only if  $w < \frac{1}{2}$ .

Consider scale similarity in previous step,  $w$  chosen with same strategy. So  $x$  should be same fraction of way from  $b - c$  as  $b$  from  $a$  to  $c$ , i.e.

$$\frac{z}{1-w} = w$$

So  $w^2 - w + (1 - 2w) = 0$  or  $w^2 - 3w + 1 = 0$  or  $w = \frac{3-\sqrt{9-4}}{2} = \frac{3-\sqrt{5}}{2} \approx 0.38197$

**33.3.3. Overall Golden Section Search algorithm.** See pp. 492 - 493, Sec. 10.2 ”Golden Section Search in One Dimension”, Press, Teukolsky, Vetterling, Flannery (2007) [36]

Given  $(a, b, c)$ , choosen new point  $x$  to be either  $(a, b)$  or  $(b, c)$ .

Then evaluate  $f(x)$ : If  $f(b) < f(x)$ , choose  $(a, b, x)$ .

If  $f(b) > f(x)$ , choose  $(b, x, c)$ .

Continue until distance between 2 outer points of triplet is tolerably small.

**33.3.4. Pseudocode for Golden Section Search.** See pp. 495 - 496, Sec. 10.2 ”Golden Section Search in One Dimension”, Press, Teukolsky, Vetterling, Flannery (2007) [36]

I will try to write with math what the code is doing for `struct Golden : Bracketmethod {}` :

Defining  $C = 0.38197$

$R = 0.61803399$

Given  $(a, b, c)$

$x_0 = a$

$x_3 = c$

if  $|c - b| > |b - a|$ ,

$x_1 = b$

$x_2 = b + C(c - b)$

else,

$x_2 = b$

$x_1 = b - C(b - a)$

while  $|x_3 - x_0| > \epsilon(|x_1| + |x_2|)$  (EY: 20230430, I don’t understand this condition and how it checks for ”tolerably small”),  
if  $f(x_2) < f(x_1)$  (then, EY: 20230430, I expected to choose  $(x_1, x_2, c)$  but the code says the following)  
 $(x_1, x_2, (1 - C)x_2 + Cx_3)$   
if  $f(x_2) \geq f(x_1)$ ,  $(x_2, x_1, (1 - C)x_1 + Cx_0)$

**33.4. Parabolic Interpolation and Brent’s Method in 1 dim.** See Sec. 10.3 ”Parabolic Interpolation and Brent’s Method in One Dimension”, Press, Teukolsky, Vetterling, Flannery (2007) [36]

On pp. 496, given Eq. (10.3.1) of Press, Teukolsky, Vetterling, Flannery (2007) [36],

(210) 
$$x = b - \frac{1}{2} \frac{(b-a)^2(f(b) - f(c)) - (b-c)^2(f(b) - f(a))}{(b-a)(f(b) - f(c)) - (b-c)(f(b) - f(a))}$$

We had derived Eq. 204. We will show that Eq. 204 and Eq. 210 are equal.  
Now Eq. 210 can be transformed into the following by factoring out a  $b$  factor from the second term of the equation:

(211) 
$$x = b - \frac{1}{2} \frac{(b-a)^2(f(b) - f(c)) - (b-c)^2(f(b) - f(a))}{(b-a)(f(b) - f(c)) - (b-c)(f(b) - f(a))} = \frac{1}{2}b - \frac{1}{2} \frac{-a(b-a)(f(b) - f(c)) + c(b-c)(f(b) - f(a))}{(b-a)(f(b) - f(c)) - (b-c)(f(b) - f(a))}$$

We will show that the denominators of Eq. 210 and Eq. 204 are equal, except for a  $(-1)$  factor:

$$(b-a)(f(b) - f(c)) - (b-c)(f(b) - f(a)) = c(f(b) - f(a)) - b(f(b) - f(a)) + b(f(b) - f(c)) - a(f(b) - f(c))$$

Now 
$$a(f(b) - f(c)) = a(f(b) - f(a) + f(a) - f(c)) = a(f(b) - f(a)) - a(f(c) - f(a))$$

So then using the very previous expression in the denominator,  
$$\implies (c-a)(f(b) - f(a)) - (b-a)(f(c) - f(a))$$

So modulo  $(-1)$ , the denominators of Eq. 210, Eq. 204, are the same!

Now 
$$a = \frac{a(b-a)(f(c) - f(a)) - a(c-a)(f(b) - f(a))}{(b-a)(f(c) - f(a)) - (f(b) - f(a))(c-a)}$$

Consider, for Eq. 204,  
$$\begin{aligned} &a(b-a)(f(c) - f(a)) - a(c-a)(f(b) - f(a)) - (c-b)(c-a)(f(b) - f(a)) = \\ &= a(b-a)(f(c) - f(a)) - (c-b+a)(c-a)(f(b) - f(a)) = \\ &= a(b-a)f(c) - a(b-a)f(a) + c(b-c)(f(b) - f(a)) - (-ac + ab + ac - a^2)(f(b) - f(a)) = \\ &= a(b-a)(f(c) - f(b)) + c(b-c)(f(b) - f(a)) \end{aligned}$$

where the last expression is because the  $-a(b-a)f(a)$  terms cancel.  
Comparing this expression on the ” $x$ ” side, Eq. 210, in Eq. 211, we have the same numerator, and since denominators are equal other than  $(-1)$  factor, so  $\frac{-1}{2}(-1) = \frac{1}{2}$ , we have equal expressions.

Consider 6 pts., not necessarily distinct,  $a, b, u, v, w, x$ .

Minimum is bracketed between  $a, b$ ;  $x$  is pt. with very least function value so far (or most recent, in case of tie).

$w \equiv$  pt. with 2nd. least function value  
 $v \equiv$  previous value of  $w$   
 $u =$  pt. which function was evaluated most recently

General principles:  
parabolic interpolation attempted, fitting through the pts.  $x, v, w$ .  
To be acceptable, parabolic step must  
(1) fall within bounding interval  $(a, b)$  and

(2) imply movement from best current value  $x$  that is *less* than half movement of *step before last*

This ensures parabolic steps actually converging to something, rather than, say, bouncing around some nonconvergent limit cycle.

Reason for comparing to step *before last* seems essentially heuristic: by experience, it’s better not to ”punish” the algorithm for a single bad step if it can make it up on next one.

34. CONJUGATE GRADIENT METHODS

See [wikipedia, Conjugate Gradient method](#).

**34.1. Problem: Conjugate Gradient Method.** Suppose we want  $\mathbf{x}$ , given  $A\mathbf{x} = \mathbf{b}$ , where  $N \times N$  matrix  $A$  is symmetric and positive definite (i.e.  $x^T Ax > 0, \forall x \neq 0$ ).

**34.2. Derivation of Conjugate Gradient Method as Direct Method. Conjugate vectors.** 2 non-zero vectors  $u, v$  are conjugate with respect to  $A$  if  $u^T Av = 0$  i.e.  $u_i A_{ij} v_j = 0$ .

Since  $A$  is symmetric and positive-definite, the left hand side defines inner product.

(212) 
$$\langle \mathbf{u}, \mathbf{v} \rangle_A \equiv \langle A\mathbf{u}, \mathbf{v} \rangle = (Au)^T v = u^T A^T v = u^T Av = \langle u, Av \rangle$$

Suppose  $P = \{p_1 \dots p_n\}$  s.t.  $p_i^T Ap_j = 0 = \langle p_i, p_j \rangle_A$  if  $i \neq j$ ; then mutually conjugate vectors  $\{p_i\}$  form a basis for  $\mathbb{R}^n$ , and so let  $x = x_*$  s.t.  $Ax_* = b$ .

$$\begin{aligned} x_* &= \sum_{i=1}^n \alpha_i p_i \implies Ax_* = \sum_{i=1}^n \alpha_i Ap_i = b \\ \implies p_k^T b &= p_k^T Ax_* = \sum_{i=1}^n \alpha_i p_k^T Ap_i = \sum_{i=1}^n \alpha_i \langle p_k, p_i \rangle_A = \alpha_k \langle p_k, p_k \rangle_A \\ &\text{or } \alpha_k = \frac{p_k^T b}{\langle p_k, p_k \rangle_A} \end{aligned}$$

So since  $x_* = \sum_{i=1}^n \alpha_i p_i$ , we have a solution:  
Find  $n$  conjugate directions,  $P$ , and then compute  $\alpha_k$ ’s.

**34.3. Conjugate Gradient as iterative.** Denote initial guess for  $x_*$  s.t.  $Ax_* = b$  to be  $x_0$  (assume  $x_0 = 0$ ; otheriwse consider  $Az = b - Ax_0$ , because if  $x_0 \neq 0$ ,  $b' := b - Ax_0$  and so for  $Az = b' = 0$ ,  $z = 0$  since  $A$  positive definite).

Because  $A$  is positive definite, consider  $Ax = 0$  and  $x \neq 0$ . But then  $x^T Ax = 0$ , contradicting  $A$  positive definite. Then  $A$  has a 0-dim. nullspace,  $\{0\}$ .

Consider 
$$f(x) = \frac{1}{2} x^T Ax - x^T b = \frac{1}{2} x_i A_{ij} x_j - x_i b_i = x^T (\frac{1}{2} Ax - b)$$

For  $x = x_*$ , s.t.  $Ax_* = b$ , then  $f(x_*) = -\frac{1}{2} x^T b$ .

Consider  $\frac{\partial^2 f(x)}{\partial x^i \partial x^j} = A_{ij}$ , since  $A$  positive definite and symmetric, then the Hessian matrix of 2nd. derivatives  $\frac{\partial^2 f(x)}{\partial x^i \partial x^j}$  is positive definite, and so  $\exists$  unique minimizer. (EY: TODO: show)

Note that 
$$\nabla f(x) = Ax - b \text{ since } \partial_i f(x) = A_{ij} x_j - b_i$$

Take  $p_0 = -\nabla f(x) = b - Ax_0$ .

Since all other  $p_j$  basis vectors will be conjugate to  $p_0$ , the gradient, hence the name conjugate gradient method.

Taking  $p_0 = b - Ax_0$ , we take the direction directly opposite the gradient.

Let residual  $r_k = b - Ax_k$  of the  $k$ th step.

Now  $r_k = -\nabla f(x_k)$ , but we want direction  $p_k$  to be conjugate to all prior directions. Thus, negative of gradient, the direction we want to go for the minimum.

Use Gram-Schmidt orthonormalization

$$(213) \quad p_k = r_k - \sum_{i < k} \frac{p_i^T A r_k}{p_i^T A p_i} p_i = r_k - \sum_{i < k} \frac{\langle p_i, r_k \rangle_A}{\langle p_i, p_i \rangle_A} p_i$$

So

$$\begin{aligned} \langle p_j, p_k \rangle &= \langle p_j, r_k \rangle - \sum_{i < k} \frac{\langle p_i, r_k \rangle_A}{\langle p_i, p_i \rangle_A} \langle p_i, p_j \rangle = \langle p_j, r_k \rangle - \frac{\langle p_j, r_k \rangle_A}{\langle p_j, p_j \rangle} = 0 \text{ if } j \neq k \\ \langle p_k, p_k \rangle &= \langle p_k, r_k \rangle - 0 = \langle p_k, r_k \rangle \end{aligned}$$

Following this direction  $p_k$

$$(214) \quad x_{k+1} = x_k + \alpha_k p_k$$

$$(215) \quad \alpha_k = \frac{p_k^T (b - Ax_k)}{p_k^T A p_k} = \frac{p_k^T r_k}{\langle p_k, p_k \rangle}$$

This is found by considering the following:

$$\begin{aligned} f(x_{k+1}) &= f(x_k + \alpha_k p_k) =: g(\alpha_k) \\ g(\alpha_k) &= \frac{1}{2} (x_k^T + \alpha_k p_k^T) A (x_k + \alpha_k p_k) - (x_k^T + \alpha_k p_k^T) b \\ g'(\alpha_k) &= 0 \text{ if } g'(\alpha_k) = \frac{1}{2} p_k^T A (x_k + \alpha_k p_k) + \frac{1}{2} (x_k^T + \alpha_k p_k^T) A p_k - p_k^T b = \\ &= \frac{1}{2} p_k^T A x_k + \frac{\alpha_k}{2} p_k^T A p_k + \frac{1}{2} x_k^T A p_k + \frac{1}{2} \alpha_k p_k^T A p_k - p_k^T b = \\ &= \alpha_k p_k^T A p_k + \frac{1}{2} \langle p_k, x_k \rangle + \frac{1}{2} \langle x_k, p_k \rangle - p_k^T b = 0 \text{ or} \\ \alpha_k \langle p_k, p_k \rangle &= p_k^T (b - Ax_k) \text{ or } \alpha_k = \frac{p_k^T (b - Ax_k)}{\langle p_k, p_k \rangle} \end{aligned}$$

Thus, by considering when  $g'(\alpha_k) = 0$ , we obtain then to follow this direction  $p_k$ :

$$(216) \quad x_{k+1} = x_k + \alpha_k p_k, \quad \alpha_k = \frac{p_k^T r_k}{\langle p_k, p_k \rangle}$$

TODO: EY (20230513) consider what this means:

$$\begin{aligned} r_i^T r_j &= (b - Ax_i)^T (b - Ax_j) = (b_k - A_{kl} x_{i,l}) (b_k - A_{km} x_{j,m}) = \\ b_k b_k - b_k A_{km} x_{j,m} - b_k A_{kl} x_{i,l} + A_{kl} x_{i,l} A_{km} x_{j,m} &= b^2 - 2b^T A x + (A x)^T A x = b^2 - 2(b - \frac{1}{2} A x)^T A x \end{aligned}$$

Let's derive the iterative step for the residuals,  $\mathbf{r}_k$ . Consider

$$r_{k+1} = b - Ax_{k+1} = b - A(x_k + \alpha_{k+1} p_{k+1}) = (b - Ax_k) - \alpha_{k+1} A p_{k+1} = r_k - \alpha_{k+1} A p_{k+1}$$

where we used the definition of a residual being the difference between  $b$  and  $Ax$ , and the iterative step for  $x_{k+1}$ .

So

$$(217) \quad r_{k+1} = r_k - \alpha_{k+1} A p_{k+1} \text{ if we had defined } x_{k+1} = x_k + \alpha_{k+1} p_{k+1}$$

But if we have  $x_{k+1}$  to be derived as in Eq. 216, then

$$(218) \quad r_{k+1} = r_k - \alpha_k A p_k$$

We will now show the following the facts. This was obtained from slides 15-16 of [EE364b, Stanford University](#).

**Theorem 2.**  $p_j^T r_k = 0 \forall j < k$  and  $r_j^T r_k = 0 \forall j < k$

*Proof.* Recall that for  $x_0 = 0$  as an initial guess without loss of generality, then

$$p_0 = b = r_0$$

$$\alpha_0 = \frac{p_0^T r_0}{\langle p_0, p_0 \rangle}$$

Let's then show the base case:

$$p_0^T r_1 = p_0^T (r_0 - \alpha_0 A p_0) = p_0^T (b - Ax_0 - \alpha_0 A p_0) = p_0^T (b - \alpha_0 A p_0) = b^2 - \alpha_0 \langle p_0, p_0 \rangle = 0$$

Suppose  $p_j^T r_k = 0 \forall j < k$ .

Consider the  $k+1$  case.

$$p_j^T r_{k+1} = p_j^T (r_k - \alpha_k A p_k) = p_j^T r_k - \alpha_k \langle p_j, p_k \rangle = 0$$

since  $p_j^T r_k = 0$  by induction base and  $\langle p_j, p_k \rangle = 0$  by definition of conjugacy.

Likewise, consider this base case:

$$r_0^T r_1 = r_0^T (r_0 - \alpha_0 A p_0) = b^2 - \alpha_0 \langle p_0, p_0 \rangle = 0$$

Assume  $r_j^T r_k = 0 \forall j < k$ .

Consider now the  $k+1$  case:

$$r_j^T r_{k+1} = r_j^T (r_k - \alpha_k A p_k) = r_j^T r_k - \alpha_k \langle r_j, p_k \rangle$$

Using Gram-Schmidt orthonormality,  $p_j = r_j - \sum_{l=0}^{j-1} \frac{\langle p_l, r_j \rangle}{\langle p_l, p_l \rangle} p_l$  so

$$\langle r_j, p_k \rangle = \langle p_j + \sum_{l=0}^{j-1} \frac{\langle p_l, r_j \rangle}{\langle p_l, p_l \rangle} p_l, p_k \rangle = 0$$

Since  $j, l < k$  and  $\langle p_l, p_k \rangle = 0$  by conjugacy of  $p_j$ 's.

Thus,  $r_j^T r_{k+1} = 0$ . Then by induction, we've shown  $r_j^T r_k = 0 \forall j < k$

□

Thus, using  $p_j^T r_k = 0$  and  $r_j^T r_k \forall j < k$ , and  $r_{k+1} = r_k - \alpha_k A p_k$ ,

$$p_k = r_k - \sum_{i=0}^{k-1} \frac{\langle p_i, r_k \rangle}{\langle p_i, p_i \rangle} p_i = r_k - \sum_{i=0}^{k-1} \frac{r_k^T \frac{(r_i - r_{i+1})}{\alpha_i}}{p_i^T \frac{(r_i - r_{i+1})}{\alpha_i}} p_i = r_k - \frac{\frac{-r_k^2}{\alpha_{k-1}} p_{k-1}}{\frac{p_{k-1}^T r_{k-1}}{\alpha_{k-1}}} = r_k + \frac{r_k^2}{r_{k-1}^T p_{k-1}} p_{k-1}$$

Also,

$$r_k^T p_k = r_k^T \left( r_k + \frac{r_k^2 p_{k-1}}{r_{k-1}^T p_{k-1}} \right) = r_k^2$$

Since  $r_k^T p_{k-1} = p_{k-1}^T r_k = 0$ .

Thus,

$$(219)$$

$$p_k = r_k + \beta_k p_{k-1} \text{ where } \beta_k = \frac{r_k^2}{r_{k-1}^2}$$

34.3.1. *Algorithm for Conjugate Gradient.* Start with a guess for  $\mathbf{x}$  for  $A\mathbf{x} = \mathbf{b}$  Without loss of generality, pick  $\mathbf{x}_0 = \mathbf{0} \in \mathbb{R}^N$ .

$$r_0 := b - Ax_0$$

If  $r_0$  sufficiently small, then return  $x_0$  as result.

$$\begin{aligned} p_0 &:= r_0 \\ k &:= 0 \end{aligned} \tag{222}$$

Repeat

$$\begin{aligned} \alpha_k &:= \frac{r_k^T r_k}{p_k^T A p_k} \\ x_{k+1} &:= x_k + \alpha_k p_k \\ r_{k+1} &= r_k - \alpha_k A p_k \end{aligned}$$

If  $r_{k+1}$  sufficiently small, then exit loop.

$$\begin{aligned} \beta_k &:= \frac{r_{k+1}^T r_{k+1}}{r_k^T r_k} \\ p_{k+1} &:= r_{k+1} + \beta_k p_k \\ k &:= k + 1 \end{aligned}$$

End repeat.

Return  $\mathbf{x}_{k+1}$  as the result.  
This is implemented also in the following manner, in [https://optimization.cbe.cornell.edu/index.php?title=Conjugate\\_gradient\\_methods](https://optimization.cbe.cornell.edu/index.php?title=Conjugate_gradient_methods) for Algorithm 3, ”Simplified version from Alg 2”, and [conjugateGradient cuda sample](#):

Given initial guess  $x_0 \in \mathbb{R}^N$

$$r_0 = b - Ax_0$$

Consider  $|r_0|^2$ . If  $|r_0|^2$  small enough, return  $x_0$  as a result.

$$\begin{aligned} p_0 &:= r_0 \\ \alpha_0 &:= \frac{r_0^2}{\langle p_0, p_0 \rangle_A} \\ x_1 &:= x_0 + \alpha_0 p_0 \\ r_1 &:= r_0 - \alpha_0 A p_0 \end{aligned} \tag{220}$$

In step  $k = 1$ ,

$$\begin{aligned} \beta_0 &:= \frac{r_1^2}{r_0^2} \\ p_1 &:= r_1 + \beta_0 p_0 \\ \alpha_1 &:= \frac{r_1^2}{\langle p_1, p_1 \rangle_A} \\ x_2 &:= x_1 + \alpha_1 p_1 \\ r_2 &:= r_1 - \alpha_1 A p_1 \end{aligned} \tag{221}$$

And so for  $k > 0$ ,

$$\begin{aligned} \beta_{k-1} &:= \frac{r_k^2}{r_{k-1}^2} \\ p_k &:= r_k + \beta_{k-1} p_{k-1} \\ \alpha_k &:= \frac{r_k^2}{\langle p_k, p_k \rangle_A} \\ x_{k+1} &:= x_k + \alpha_k p_k \\ r_{k+1} &:= r_k - \alpha_k A p_k \end{aligned}$$

Repeat until  $r_k^2$  small enough or when  $k > K$ ,  $K \equiv$  maximum number of iterations.

34.4. **Biconjugate gradient method; Biconjugate gradient stabilized method.** See [Math 5330 Computational Methods of Linear Algebra, Spring 2017](#).

In [Lecture Note 6, 02/09 Lecture](#) of Math 5330, in considering a transpose free version of the BCG (biconjugate algorithm), notice  $\exists$  polynomials of degree  $j$ ,  $\phi_j$ ,  $\pi_j$  s.t.  $\phi_j(0) = 1$  and

$$\begin{aligned} \mathbf{r}_j &= \phi_j(A) \mathbf{r}_0 & \mathbf{p}_j &= \pi_j(A) \mathbf{r}_0 \\ \mathbf{r}'_j &= \phi_j(A^T) \mathbf{r}_0 & \mathbf{p}'_j &= \pi_j(A^T) \mathbf{r}_0 \end{aligned} \tag{223}$$

Eq. [223](#) is a reasonable generalization. Recall Eq. [220](#), [221](#), [222](#):

$$\begin{aligned} \mathbf{p}_0 &:= \mathbf{r}_0 \\ \mathbf{r}_1 &:= \mathbf{r}_0 - \alpha_0 A \mathbf{p}_0 = (1 - \alpha_0 A) \mathbf{r}_0 \\ \mathbf{p}_1 &:= \mathbf{r}_1 + \beta_0 \mathbf{p}_0 = ((1 - \alpha_0 A) + \beta_0) \mathbf{p}_0 \\ \mathbf{p}_k &:= \mathbf{r}_k + \beta_{k-1} \mathbf{p}_{k-1} \\ \mathbf{r}_{k+1} &:= \mathbf{r}_k - \alpha_k A \mathbf{p}_k \end{aligned}$$

So by induction

$$\begin{aligned} \mathbf{r}_{k+1} &= \phi_k(A) \mathbf{r}_0 - \alpha_k A \pi_k(A) \mathbf{r}_0 = (\phi_k(A) - \alpha_k A \pi_k(A)) \mathbf{r}_0 =: \phi_{k+1}(A) \mathbf{r}_0 \\ \mathbf{p}_k &= \phi_k(A) \mathbf{r}_0 + \beta_{k-1} \pi_{k-1}(A) \mathbf{r}_0 = (\phi_k(A) + \beta_{k-1} \pi_{k-1}(A)) \mathbf{r}_0 =: \phi_k(A) \mathbf{r}_0 \end{aligned}$$

So we get

$$\begin{aligned} \phi_{j+1}(A) &= \phi_j(A) - \alpha_j A \pi_j(A) \\ \pi_j(A) &= \phi_j(A) + \beta_{j-1} \pi_{j-1}(A) \end{aligned} \tag{224}$$

[Wikipedia](#) defines  $\mathbf{r}_i = P_i(A) \mathbf{r}_0$  and  $\mathbf{p}_{i+1} = T_i(A) \mathbf{r}_0$ , and so  $\phi_j \equiv P_j$ , and  $\pi_i \equiv T_i$ . We’ll use Wikipedia’s notation from here. We wish to have a recurrence relationship as such:

$$\begin{aligned} \mathbf{r}_j &= Q_j(A) P_j(A) \mathbf{r}_0 \\ \mathbf{p}_j &= Q_j(A) T_j(A) \mathbf{r}_0 \end{aligned}$$

$Q_j(A)$  is a called a ”rest polynomial” by Xianyi Zeng in [Lecture Note 7, Math 5330, UTEP](#). Wikipedia says the hope or wish is that  $Q_i(A)$  will enable faster and smoother convergence in  $\mathbf{r}_i$ . Consider  $Q_j(A)$  to be an ”ansatz” or guess. We already know that  $Q_j(A) = 1$  gives us back the ”original” conjugate gradient method.

From  $Q_{j+1}(A) = (1 - \omega_j A) Q_j(A)$ , Eq. [224](#),

$$\begin{aligned} \mathbf{r}_{j+1} &= Q_{j+1}(A) P_{j+1}(A) \mathbf{r}_0 = (1 - \omega_j A) Q_j(A) P_{j+1}(A) \mathbf{r}_0 = \sum_{k=0}^{j+1} a_k^{j+1} A^{j+1-k} P_{j+1}(A) \mathbf{r}_0 = (1 - \omega_j A) \sum_{k=0}^j a_k^j A^{j-k} P_{j+1}(A) \mathbf{r}_0 \\ & \xrightarrow{\mathbf{r}_0^T} a_0^{j+1} = -\omega_j a_0^j \end{aligned} \tag{225}$$

We’ll compare a number of implementations of the algorithm for the biconjugate gradient stabilized (BCGST). Consider this implementation from Zeng, pp.5 Algorithm 2.3, in [Lecture Note 7, Math 5330, UTEP](#).

(226)

Compute  $\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0$ , choose  $\mathbf{r}'_0$  s.t.  $\mathbf{r}_0 \cdot \mathbf{r}'_0 \neq 0$ , e.g.  $\mathbf{r}'_0 = \mathbf{r}_0$   
Set  $\mathbf{p}_0 = \mathbf{r}_0$

$$\begin{aligned} &\forall j = 0, 1, \dots \\ &\alpha_j = \frac{(\mathbf{r}'_0 \cdot \mathbf{r}_j)}{\mathbf{r}'_0 \cdot A\mathbf{p}_j} \\ &\mathbf{s}_j = \mathbf{r}_j - \alpha_j A\mathbf{p}_j \\ &\omega_j = \frac{\mathbf{s}_j \cdot (A\mathbf{s}_j)}{(A\mathbf{s}_j) \cdot (A\mathbf{s}_j)} \\ &\mathbf{x}_{j+1} = \mathbf{x}_j + \alpha_j \mathbf{p}_j + \omega_j \mathbf{s}_j \\ &\mathbf{r}_{j+1} = \mathbf{s}_j - \omega_j A\mathbf{s}_j \\ &\text{if } \|\mathbf{r}_{j+1}\| < \epsilon_0 \text{ then break.} \\ &\beta_j = \frac{\alpha_j}{\omega_j} \frac{(\mathbf{r}'_0 \cdot \mathbf{r}_{j+1})}{(\mathbf{r}'_0 \cdot \mathbf{r}_j)} \\ &\mathbf{p}_{j+1} = \mathbf{r}_{j+1} + \beta_j (\mathbf{p}_j - \omega_j A\mathbf{p}_j) \\ &\text{Set } \mathbf{x} = \mathbf{x}_{j+1} \end{aligned}$$

Compute  $\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0$ , choose  $\mathbf{r}'_0$  s.t.  $\mathbf{r}_0 \cdot \mathbf{r}'_0 \neq 0$ , e.g.  $\mathbf{r}'_0 = \mathbf{r}_0$   
Set  $\mathbf{p}_0 = \mathbf{r}_0$

(227)

$$\begin{aligned} &\forall j = 0, 1, \dots \\ &\alpha_j = \frac{(\mathbf{r}'_0 \cdot \mathbf{r}_j)}{\mathbf{r}'_0 \cdot A\mathbf{p}_j} \\ &\mathbf{s}_j = \mathbf{r}_j - \alpha_j A\mathbf{p}_j \\ &\text{if } \|\mathbf{s}_j\| < \epsilon_0 \text{ , then} \\ &\quad \mathbf{x}_{j+1} = \mathbf{x}_j + \alpha_j \mathbf{p}_j \\ &\quad \text{break} \\ &\omega_j = \frac{\mathbf{s}_j \cdot (A\mathbf{s}_j)}{(A\mathbf{s}_j) \cdot (A\mathbf{s}_j)} \\ &\mathbf{x}_{j+1} = \mathbf{x}_j + \alpha_j \mathbf{p}_j + \omega_j \mathbf{s}_j \\ &\mathbf{r}_{j+1} = \mathbf{s}_j - \omega_j A\mathbf{s}_j \\ &\text{if } \|\mathbf{r}_{j+1}\| < \epsilon_0 \text{ then break.} \\ &\beta_j = \frac{\alpha_j}{\omega_j} \frac{(\mathbf{r}'_0 \cdot \mathbf{r}_{j+1})}{(\mathbf{r}'_0 \cdot \mathbf{r}_j)} \\ &\mathbf{p}_{j+1} = \mathbf{r}_{j+1} + \beta_j (\mathbf{p}_j - \omega_j A\mathbf{p}_j) \\ &\text{if } |\mathbf{r}_{j+1} \cdot \mathbf{r}'_0| < 10^{-6} \text{ then} \\ &\quad \mathbf{r}'_0 = \mathbf{r}_{j+1} \\ &\quad \mathbf{p}_{j+1} = \mathbf{r}_{j+1} \\ &\quad \text{Set } \mathbf{x} = \mathbf{x}_{j+1} \end{aligned}$$

We will use this algorithm for the biconjugate gradient stabilized (BCGST)

Compute  $\mathbf{r}_0 = \mathbf{b} - A\mathbf{x}_0$ , choose  $\mathbf{r}'_0$  s.t.  $\mathbf{r}_0 \cdot \mathbf{r}'_0 \neq 0$ , e.g.  $\mathbf{r}'_0 = \mathbf{r}_0$   
Set  $\mathbf{p}_0 = \mathbf{r}_0$

(228)

$$\begin{aligned} &\forall j = 0, 1, \dots \\ &\alpha_j = \frac{(\mathbf{r}'_0 \cdot \mathbf{r}_j)}{\mathbf{r}_0 \cdot A\mathbf{p}_j} \\ &\mathbf{s}_j = \mathbf{r}_j - \alpha_j A\mathbf{p}_j \end{aligned}$$

Part 9. Cantera

35. PROFESSOR DAVID GOODWIN, 1957-2012

“Keep on working to make this world a better place. It is important to get your thoughts away from your problems and focused on those who are less fortunate than you are.”  
“We can’t control the hand of cards we are dealt, but we sure can determine how we play them.”

- Open-Source, Berkeley license,
- Object-Oriented Structure, multi-inferface structure,
- formalism for collision integrals in evaluating transport parameters of dilute ideal gases

cf. <https://code.google.com/archive/p/cantera/wikis/DaveGoodwin.wiki>

The following includes (disorganized) notes on using, implementing, and developing for Cantera.

Unfortunately, this algorithm, Algorithm 2.3, will fail and yield NaN values especially during floating point division. Use this instead, Algorithm 2.4 of Zeng:



EY : 20160125 I wanted to obtain the *equation of state* but in Python. This was implemented in Matlab as `eosType.m` here <https://github.com/Cantera/cantera/blob/master/interfaces/matlab/toolbox/%40ThermoPhase/eosType.m> It leads to `thermo_get` for the code says

```
e = thermo_get(tp.tp_id , 18);
```

I don’t think the `tp.tp_id` equivalent in Python is `tp.ID`; searching for `tp_id` led to [ThermoPhase.m](#) and

```
if nargin == 1
    if isa(r, 'ThermoPhase')
        % create a copy
        t = r;
        return
    elseif isa(r, 'XMLNode')
        t.owner = 1;
        hr = xml_hndl(r);
        t.tp_id = thermo_get(hr, 0);
        if t.tp_id < 0
            error(geterr);
        end
    else
        t.owner = 0;
        t.tp_id = r;
    end
    t = class(t, 'ThermoPhase');
```

`thermo_get` of MatLab is here, in the `private` folder: [https://github.com/Cantera/cantera/blob/master/interfaces/matlab/toolbox/%40ThermoPhase/private/thermo\\_get.m](https://github.com/Cantera/cantera/blob/master/interfaces/matlab/toolbox/%40ThermoPhase/private/thermo_get.m) which leads to the MatLab command `ctmethods` and the following inputs/parameters

```
function i = thermo_get(n, job, a, b)
if nargin == 2
    i = ctmethods(20, n, job);
```

I was only able to find the “definition” of `ctmethods` for *MatLab* here: <https://github.com/Cantera/cantera/blob/5926d2db7c0d4919b75ee50828b0adab4e691a51/src/matlab/ctmethods.cpp>

which was essentially a so-called “mex function” which was a bunch of different cases; the case I was interested in was for Thermo, leading me to <https://github.com/Cantera/cantera/blob/5926d2db7c0d4919b75ee50828b0adab4e691a51/src/matlab/thermomethods.cpp> with its `thermo_get` and `thermo_set`

as above, the so-called “job” I wanted was `job=18`, which lead me to the command `th_eosType`

```
case 18:
    vv = double(th_eosType(n));
    break;
```

In [soundspeed.m](#) for the Matlab toolbox ThermoPhase, notice the line `if isIdealGas(tp)`. This is implemented in [isIdealGas.m](#) by returning `v` in

```
if eosType(tp) == 1
    v = 1;
else
    v = 0;
```

### 36. SPEED OF SOUND AND SOUND SPEED.PY

Take a look at `sound_speed.py` in the github repository for [cantera](#).  
Take a look at `afrozen` in the code. This is an exact implementation of the usual speed of sound:

$$a_{\text{frozen}}^2 = \left(\frac{\partial p}{\partial \rho}\right)_\sigma$$

Take a look at `afrozen2` in the code. This is an exact implementation of the speed of sound for an ideal gas. Let’s review its derivation.

From  $pV = N\tau$ , which is how ideal gases behave at any time, then  $p = \frac{\rho\tau}{M}$  where  $\rho := \frac{MN}{V}$  and  $M$  is the mass of a single particle of the species. Then

$$dp = \frac{\tau d\rho + \rho d\tau}{M}$$

and so

$$C_p \frac{d\tau}{\tau} = C_p \left(\frac{dp}{p} - \frac{d\rho}{\rho}\right) = \frac{V dp}{\tau}$$

where the last equality comes from the following derivation

$$\begin{aligned} Q &= C_p d\tau - V dp = \tau d\sigma = 0 \text{ since } d\sigma = 0 \text{ for an isentropic process} \\ \implies \frac{C_p}{\tau} d\tau &= \frac{V}{\tau} dp \end{aligned}$$

Thus

$$C_p \frac{d\rho}{\rho} = \left(\frac{C_p}{p} - \frac{V}{\tau}\right) dp = \left(\frac{C_p - N}{p}\right) dp = C_V \frac{dp}{p} \implies \left(\frac{\partial p}{\partial \rho}\right)_\sigma = \frac{\gamma\tau}{M}$$

recovering the usual speed of sound for ideal gases.

The derivation on pp. 3-4 of Lecture 14 of MIT OCW 16.512 by Martinez-Sanchez (2005) [20] is fallacious because  $M$  is constant always and even in chemical reactions, masses for each species don’t get created or destroyed. The changes due to chemical reactions show up in  $dN_i$ . And even then, we know that the changes in  $N_i$ , the number of particles for each species  $i$ ,  $dN_i$ , is governed by the chemical reaction and the number of times the chemical reaction occurs.

For instance, surely the total mass  $m$  of the system is thus given by

$$m \equiv \sum_j m_j := \sum_j N_j M_j$$

and so

$$dm = \sum_j dm_j = \sum_j M_j dN_j = \sum_j M_j \nu_j d\hat{N}$$

where  $d\hat{N} \in \mathbb{Z}$  indicates how many times a reaction occurs, and  $dN_j = \nu_j d\hat{N}$ , which tells us how the number of particles in each species changes everytime a chemical reaction occurs, as the chemical reaction is given by  $\sum_j \nu_j A_j = 0$  (see pp. 267 Section **Equilibrium in Reactions**, Ch. 9 Gibbs Free Energy and Chemical Reactions of Kittel and Kroemer [8]).

Considering the total number of particles  $N = \sum_j N_j$ , the total pressure  $p = \sum_j p_j$  as a sum of partial pressure  $p_j$ , and taking the system as a whole to always obey ideal gas behavior,  $pV = N\tau$ , then

$$p_j V = N_j \tau$$

so that each species, at any time, behaves as an ideal gas.

Considering the (mass) density of species  $i$ ,  $\rho_j := \frac{N_j M_j}{V}$ , and so  $\sum_j \rho_j = \frac{m}{V} = \rho$ , then rewrite the ideal gas law as  $p_j = \frac{N_j \tau}{V} = \frac{\rho_j \tau}{M_j}$ .

If one takes the differential as does Martinez-Sanchez (2005) [20] in Lecture 14, pp. 3-4,

$$dp_j = \frac{\tau d\rho_j + \rho_j d\tau}{M_j} \implies \frac{dp_j}{p_j} = \frac{d\rho_j}{\rho_j} + \frac{d\tau}{\tau}$$

and the derivation for speed of sound  $a^2 = \left(\frac{\partial p_j}{\partial \rho_j}\right)_\sigma$  goes as above. It is fallacious to say that  $M_j$  changes as it is a constant for each species  $j$ . The changes in  $N_j$  due to chemical reactions are already accounted for in  $\rho_j$ .

Instead, one must go back to the general form of heat  $Q$ :

$$Q = Q(\tau, p, \{N_i\}) = C_p d\tau + \left(\frac{\partial \tau \sigma}{\partial p}\right)_{\tau, \{N_i\}} + \sum_i \frac{\partial Q}{\partial N_i} dN_i = \tau d\sigma$$

and realize that the implementation of `aequil` in `sound_speed.py` is this general implementation of

$$a^2 = \left( \frac{\partial p}{\partial \rho} \right)_{\sigma, \{N_i\}}$$

that does not assume ideal gas behavior, and also equilibrates the changes in species,  $dN_i$ , and does not presume that heat capacities  $C_V, C_p$  have the nice form that arises from the assumption of ideal gas behavior.

## Part 10. Compressible Flow

I'll review or verify aspects of *compressible flow* for gas dynamics.

I want to verify the so-called *Bernoulli constant* or Bernoulli invariant  $h + u^2/2 + \Phi$  on pp. 38 of Thompson (1988) [19] for *compressible* flow. Note that Thompson uses  $h := H/m$  to denote the *specific enthalpy*, whereas I used the notation of  $h := H/V$  for the enthalpy density. I'll use the same notation and one should be able to understand which one is being used from context and/or from the units.

Now for the (internal) energy density  $\epsilon$  and kinetic energy density  $k$ , defined such that

$$\begin{aligned} \epsilon &:= \frac{E}{V} \\ k &= \frac{1}{2} \rho u^2 \end{aligned} \implies E_{\text{tot}} = \int_V (\epsilon + k) \text{vol}^n$$

(which is the notation used by Le Bellac, Mortessagne, Batrouni (2004) [15]), then from the physical principle that the change in total energy of a system is given by the work done on it by the sum of all external forces on the system,

$$\frac{d}{dt} \int_V (\epsilon + k) \text{vol}^n = \int_V \left[ \frac{\partial}{\partial t} (\epsilon + k) + \text{div}((\epsilon + k)\mathbf{u}) \right] \text{vol}^n = \int_V \rho \langle u, b \rangle \text{vol}^n + \int_{\partial V} \langle u, T^j \rangle dS_j - \int_{\partial V} q^j dS_j$$

where  $b$  is the density associated with (external) body forces on the system and  $q$  is the outward-directed heat flux representing the system's heat loss, i.e.

$$q^j dS_j = q \cdot n > 0 \text{ for outward-directed heat flux,} \quad - \int_{\partial V} q^j dS_j \text{ represents heat loss}$$

The *specific energy* (unfortunately denoted as  $e$  in Thompson (1988) [19], as it looks like the exponential), the energy per unit mass, or specific internal energy, a thermodynamic quantity, is defined as such:

$$\rho e = \epsilon = \frac{MN}{V} e = \frac{E}{V} \implies e := \frac{E}{m}$$

Then

$$\frac{d}{dt} \int_V (\rho e + k) \text{vol}^n = \int_V \left[ \frac{\partial}{\partial t} (\rho e + k) + \text{div}((\rho e + k)\mathbf{u}) \right] \text{vol}^n \implies \frac{D}{Dt} (\rho e + k) + (\rho e + k) \text{div} u = \text{div} \langle u, T \rangle + \rho \langle u, b \rangle - \text{div} q$$

If one can use mass conservation for the flow

$$\frac{\partial \rho}{\partial t} + \text{div}(\rho u) = \frac{D\rho}{Dt} + \rho \text{div} u = 0 \text{ for } \dot{m} = 0$$

then

$$\rho \frac{D}{Dt} \left( e + \frac{k}{\rho} \right) = \text{div} \langle u, T \rangle + \rho \langle u, b \rangle - \text{div} q$$

since, from Newton's second law and using that on  $\frac{Du^2}{Dt}$ ,

$$\begin{aligned} \rho \frac{Du}{Dt} &= \text{div}(T) \\ \rho \frac{Du^2}{Dt} &= \langle u, \text{div} T \rangle \end{aligned}$$

Working the following out,

$$\text{div} \langle u, T \rangle = \frac{1}{\sqrt{g}} \frac{\partial}{\partial x^k} (g_{ij} u^i T^{jk} \sqrt{g}) = \langle u, \text{div} T \rangle + T^{jk} \frac{\partial}{\partial x^k} (g_{ij} u^i)$$

For  $T = -pg$ ,

$$\text{div} \langle u, T \rangle = \text{div}(g_{ij} u^i T^j) = \frac{1}{\sqrt{g}} \frac{\partial (g_{ij} u^i T^{jk} \sqrt{g})}{\partial x^k} = \frac{1}{\sqrt{g}} \frac{\partial (g_{ij} u^i (-p) g^{jk} \sqrt{g})}{\partial x^k} = \frac{-1}{\sqrt{g}} \frac{\partial (u^k p \sqrt{g})}{\partial x^k} = -u^j \frac{\partial p}{\partial x^j} - p \text{div} u$$

and

$$\frac{D}{Dt} \left( \frac{p}{\rho} \right) = \frac{1}{\rho} \frac{D}{Dt} p + \frac{-p}{\rho^2} \frac{D}{Dt} (\rho) = \frac{1}{\rho} \frac{Dp}{Dt} - \frac{p}{\rho} \text{div} u$$

since, again, using mass conservation,

$$\frac{\partial \rho}{\partial t} + \text{div}(\rho u) = \frac{D\rho}{Dt} + \rho \text{div} u = 0$$

Then

$$\rho \frac{D}{Dt} \left( \frac{p}{\rho} \right) = \frac{Dp}{Dt} - p \text{div} u$$

and so

$$\implies \rho \frac{D}{Dt} \left( e + \frac{p}{\rho} + \frac{k}{\rho} \right) = \frac{\partial p}{\partial t} + -q$$

in the case of  $b = 0$  (no external body forces).

### 37. NOZZLE FLOW OF A REACTING GAS

This follows from Section 3.7 **Nozzle Flow of a Reacting Gas** of Oates (1997) [1]. Also see Lecture 14 of Martinez-Sanchez (2005) [20].

Investigate 2 limits to flow within nozzles, “equilibrium” and “frozen.”

**37.1. Equilibrium flow.** Suppose all chemical reactions occur in times very short compared to the time of fluid passage through the nozzle, so that the fluid will be at all times (almost) in a state of chemical equilibrium. Clearly, this is an approximation as we can imagine that during supersonic flow, the reactants could've traveled far, already, before the chemical reaction goes to completion.

The reactions occur continuously throughout the flow, leading to a continuous passage of energy from chemical binding and excitation modes to translational modes. Since the fluid is kept at equilibrium (locally) at all times, the “equivalent” temperatures of all such modes of energy storage are equal and as a result the total entropy of the fluid remains constant.” (cf. pp. 86 of Oates (1997) [1]).

From Oates (1997) [1], starting from

$$\begin{aligned} \tau d\sigma_j &= C_{p_j} d\tau - V dp_j \text{ or } d\sigma_j = \frac{C_{p_j} d\tau}{\tau} - \frac{N_j dp_j}{p_j} \\ \implies \tau d\sigma_j &= C_{p_j} d\tau - \frac{N_j \tau}{p_j} dp_j \text{ or } \hat{\sigma}_j - \hat{\sigma}_j(p_0, \tau_0) = \int_{\tau_0}^{\tau} \frac{\hat{C}_{p_j} d\tau}{\tau} - \ln \frac{p_j}{p_{j0}} \\ \hat{\sigma}_j - \hat{\sigma}_j(p_0, \tau_0) &= \int_{\tau_0}^{\tau} \frac{\hat{C}_{p_j} d\tau}{\tau} - \ln \frac{p_j}{p_{j0}} \\ \sigma &= \sum_{j=1}^{\mathcal{N}} N_j \hat{\sigma}_j = \sum_{j=1}^{\mathcal{N}} N_j \left[ \left( \int_{\tau_0}^{\tau} \frac{\hat{C}_{p_j} d\tau}{\tau} + \hat{\sigma}_j(p_0, \tau_0) \right) - \ln \frac{p_j}{p_{j0}} \right] \end{aligned}$$

Recall that a chemical reaction, labeled by  $I$ , with the convention that reactant (stoichiometric) coefficients are negative integers, is given by  $\sum_j \nu_{jI} A_j = 0$  and so

$$\sum_j \nu_{jI} A_j = 0 \implies dN_{jI} = \nu_{jI} d\hat{N}_I$$

which informs us that a chemical reaction, every time it occurs, adds in and subtracts out particles of each species, represented by  $dN_{jI}$ , each time a chemical reaction is run,  $d\hat{N}$ .

Recalling the definition of enthalpy  $H = H(\sigma, p; \{N_i\}_i) := U + pV$ , then the total enthalpy  $H$  of a system is completely specified by total entropy  $\sigma$ , pressure  $p$ , and  $\mathcal{N}$  numbers,  $\{N_i\}_i \equiv \{N_i\}_{i=1\dots\mathcal{N}}$  for the number of particles in each species, with  $\mathcal{N}$  being the total number of (different) species to consider in all possible chemical reactions  $I$ . Effectively,  $H$  is a smooth scalar function on a  $2 + \mathcal{N}$  manifold  $\Sigma$ , with  $\Sigma$  representing all possible thermodynamic processes. Also, recall that

$$dH = \tau d\sigma + V dp + \sum_i \mu_i dN_i$$

Consider a curve on  $\Sigma$ ,  $\gamma : \mathbb{R} \rightarrow \Sigma$ , representing a thermodynamic process such that  $\sigma$  is constant (isentropic) and  $p$  is constant (isobaric). This means that  $d\sigma(\dot{\gamma}) = 0$  and  $dp(\dot{\gamma}) = 0$ , respectively, and so we have  $d\sigma = dp = 0$  in this case. Then, this case,

$$\begin{aligned} dH &= \sum_i \mu_i dN_i \\ \implies dH(\dot{\gamma}) &= \sum_i \mu_i \Delta N_i = H_P(\tau_P) - H_R(\tau_0) \end{aligned}$$

where  $H_P, H_R$  refer to the total enthalpies of the products and reactants, respectively, as defined in Eq. 55.

Recall that  $dG = dH - d(\tau\sigma)$  and  $dF = dU - d(\tau\sigma) = 0$ . Suppose  $dH = 0$  and  $dG = 0$ . This implies that  $d(\tau\sigma) = 0$  so therefore  $dU = 0$  and  $dF = 0$  and so we're in Helmholtz free energy equilibrium. One can imagine that both  $\sigma$  and  $\tau$  change during the process in such a way that  $d(\tau\sigma) = 0$ .

If we suppose we want Gibbs equilibrium  $dG = 0$ , and so for an isentropic process  $d\sigma = 0$ , then  $dH = \sigma d\tau$  and so the so-called *stagnation* enthalpy can drop (or at least change) with a drop (or change) in temperature.

Thus, what's going on is that there's a transformation, or mapping, via a thermodynamic process represented by curve  $\gamma : \mathbb{R} \rightarrow \Sigma$  in  $\Sigma$ , a  $2 + \mathcal{N}$ -dim. manifold:

$$\begin{array}{ccc} \Sigma & \xrightarrow{\gamma} & \Sigma \\ & \gamma \text{ s.t.} & \\ & dG(\dot{\gamma})=0 & \\ & d\sigma(\dot{\gamma})=0 & \\ & dp(\dot{\gamma})=0 & \end{array} \quad (\tau_0, p_0, \{N_{i0}\}) \mapsto (\tau_p, p_0, \{N_{ieq}\})$$

and if we wanted enthalpy to remain constant,

$$\begin{array}{ccc} \Sigma & \xrightarrow{\gamma} & \Sigma \\ & \gamma \text{ s.t.} & \\ & dG(\dot{\gamma})=0 & \\ & dH(\dot{\gamma})=0 & \\ & d(\tau\sigma)(\dot{\gamma})=0 & \end{array} \quad (\tau_0, p_0, \{N_{i0}\}) \mapsto (\tau_p, p_0, \{N_{ieq}\})$$

I think what we want, since the flow was shown to be isentropic, and, again, as Oates [1] argued, that local equilibrium leads to equilibrium in all modes of energy storage, chemical bonds, vibrational modes, to translational modes, that the total entropy of

the fluid remains constant, this:

$$\begin{array}{ccc} \Sigma & \xrightarrow{\gamma} & \Sigma \\ & \gamma \text{ s.t.} & \\ & dG(\dot{\gamma})=0 & \\ & d\sigma(\dot{\gamma})=0 & \\ & dp(\dot{\gamma})=0 & \end{array} \quad (\tau_0, p_0, \{N_{i0}\}) \mapsto (\tau_p, p_0, \{N_{ieq}\})$$

37.1.1. *Continuity equations for equilibrium flow.* The continuity equations to use for equilibrium flow should be (correct me if I'm wrong) are

$$\dot{m} = \rho u A$$

for which  $\dot{m}$  should be constant at each point along the 1-dimensional flow (out) by mass conservation (which should still should hold for chemical reactions).

One possible form of mass conservation that could prove useful is

$$\rho_2 u_2 A_2 = \rho_1 u_1 A_1 \text{ or } \frac{A_2}{A_1} = \frac{\rho_1 u_1}{\rho_2 u_2}$$

Also, throughout the flow, the Bernoulli invariant is  $h + k$ , i.e.

$$\frac{u_2^2}{2} + h_2 = \frac{u_1^2}{2} + h_1 \text{ or } u_2^2 = u_1^2 + 2(h_1 - h_2)$$

Note that the speed of sound  $a^2 = \left(\frac{\partial \rho}{\partial p}\right)_{\sigma, \{N_i\}}$  can be calculated at each instance, and so the Mach number can be calculated:

$$\mathfrak{M} = \frac{u}{a}$$

37.2. **Frozen flow.** Note that there are a number of critical typos in Martinez-Sanchez (2005) [20] that makes it difficult to understand, in this case, Lecture 14.

Given nozzle exit (i.e. entering combustion chamber), pressure  $p_c$ , and consider chamber entropy  $\sigma_c$ . Then  $(\sigma_c, p_c)$  specifies completely the thermodynamic state, and so enthalpy (or specific enthalpy)  $h_c$  is specified.

$$(\sigma_c, p_c) \mapsto h_c(\sigma_c, p_c)$$

For adiabatic flows, the energy equation (for steady state) is

$$\frac{u_2^2}{2} = h_c - h_2$$

In general, for multiple number of species,

$$(\sigma_c, p_c; \{N_i\}_{i=1\dots\mathcal{N}}) \mapsto h_c(\sigma_c, p_c; \{N_i\})$$

For throat conditions, consider mass flow  $\dot{m} = \rho u A$ .

Assuming steady state (EY: 20160210 check if throat conditions necessitates  $\dot{m} = 0$ ), then

$$\begin{aligned} \frac{d\rho}{\rho} + \frac{du}{u} + \frac{dA}{A} &= 0 \implies \frac{du}{u} = -\frac{dh}{2(h_c - h)} = \frac{-dp}{2\rho(h_c - h)} = \frac{-dp}{\rho u^2} \text{ since} \\ u du &= -dh \text{ and} \end{aligned}$$

$$dh = \tau ds + \frac{1}{\rho} dp + \frac{h}{\rho} d\rho \text{ and for this thermodynamic process } ds = d\rho = 0 \text{ in the fluid-at-rest frame}$$

At the throat,  $dA = 0$  (by definition of a throat), and so

$$\frac{d\rho}{\rho} = -\frac{du}{u} = \frac{dp}{\rho u^2} \implies (u^*)^2 = \left(\frac{\partial p}{\partial \rho}\right)_\sigma$$

Also, consider from mass conservation, that, at the throat,  $dA = 0$ ,

$$\begin{aligned} d(\rho u)A + \rho u dA &= 0 \\ \implies \frac{d(\rho u)}{\rho u} + \frac{dA}{A} &= 0 \xrightarrow{dA=0} d(\rho u) = 0 \end{aligned}$$

$\rho u$  is at a *maximum* at the throat.

**37.3. Implementation in Cantera of equilibrium flow and frozen flow.** The first thing that happens is *combustion*! Fuel and oxidizer enters through inlets into combustion chamber. Then combustion occurs in the (combustion) chamber!

$$\begin{aligned} \Sigma &\xrightarrow{\gamma} \Sigma \\ (\tau_0, p_0, \{N_{i0}\}) &\xrightarrow[\substack{\gamma \text{ s.t.} \\ dH(\dot{\gamma})=0 \\ dp(\dot{\gamma})=0}]{\gamma} (\tau_c, p_0, \{N_{ic}\}) \end{aligned}$$

and  $dH(\dot{\gamma}) = 0$  since  $H_P(\tau_c) = H_R(\tau_0)$  so  $\Delta H \equiv H_p(\tau_c) - H_R(\tau_0) = 0$   
 $h(\sigma_c, p_0, \{N_{ic}\}_i) \equiv h_0 \in C^\infty(\Sigma)$  and  $h_0$  is the stagnation enthalpy, used as the Bernoulli invariant in the energy equation (i.e.  $\frac{u^2}{2} = (h_0 - h)$ ).

Then consider isentropic flow  $\gamma_{\text{isent}}$ :

$$\begin{aligned} \Sigma &\xrightarrow{\gamma_{\text{isent}}} \Sigma \\ (\tau_c, p_0, \{N_{ic}\}) &\xrightarrow[\substack{\gamma_{\text{isent}} \text{ s.t.} \\ d\sigma(\dot{\gamma}_{\text{isent}})=0}]{\gamma_{\text{isent}}} (\tau, p, \{N_{if}\}) \\ C^\infty(\Sigma) &\xrightarrow{\gamma_{\text{isent}}} C^\infty(\Sigma) \\ h(\sigma_c, p_0, \{N_{ic}\}) \equiv h_0 &\xrightarrow{\gamma_{\text{isent}}} h(\sigma_c, p, \{N_{if}\}) \end{aligned}$$

How I implemented this in Cantera is by setting the state of the gas to the desired  $p$ , and then equilibrating the gas under constant entropy  $\sigma_c$  and constant pressure  $p$ .

**37.3.1. Getting the results for rocket (and nozzle) performance (i.e. collecting the dividends).** You want to calculate the specific impulse (for the rocket). Use this handy relation:

$$gI_{\text{sp}} = \frac{F}{\dot{m}} = \frac{\dot{m}u_{\text{exh}} + (p_{\text{exh}} - p_a)A_e}{\rho_{\text{exh}}u_{\text{exh}}A_e} = u_{\text{exh}} + \frac{p_{\text{exh}} - p_a}{\rho_{\text{exh}}u_{\text{exh}}}$$

and  $\dot{m} = \rho^* u^* A^*$ .

38. DYNAMIC PRESSURE

cf. **NASA Glenn Research Center, "Dynamic Pressure"**

From the NASA Glenn Research Center website link, here’s how they started their derivation:  
By fluid momentum conservation:

(229) 
$$\frac{-dp}{dx} = \rho u \frac{du}{dx}$$

where  
 $p \equiv$  pressure

$\rho \equiv$  density.

If  $u \in \mathbb{R}$ ,

$$\frac{dp}{dx} + \rho u \frac{du}{dx} = \frac{d}{dx} (p + \frac{1}{2} \rho u^2) = 0 \implies p_s + \frac{1}{2} \rho u^2 = \text{constant} = p_{\text{tot}}$$

where  
 $p_{\text{tot}} :=$  total pressure  
 $p_s :=$  static pressure  
 $q :=$  dynamic pressure  $= \frac{1}{2} \rho u^2$ .

If a gas is static and not flowing the measured pressure is the same in all directions.  
If gas is moving, measured pressure depends on direction of motion.

(230) 
$$p_s + \frac{1}{2} \rho u^2 = p_{\text{tot},0}$$

*looks like* incompressible Bernoulli’s equation.  
 $\frac{1}{2} \rho u^2$  is called dynamic pressure because it’s a pressure term associated with velocity of the flow.

Dynamic pressure is a defined property of moving flow of gas.  
We can use and apply idea of dynamic pressure in much more complex flows, like compressible flows or viscous flows. Particularly, aerodynamic forces acting on object as it moves through air are directly proportional to dynamic pressure.  
By measuring dynmaic pressure in flight, pitot-static tube (Prandtl tube) can be used to determine aircraft airspeed.

My question is this: how does this derivation change if we can’t assume flow is uniform and we can’t assume flow is all in 1-dimension (i.e. it has non-zero velocity components in the ”other” direction)?

Recall Navier-Stokes for compressible, viscous fluid flow:

$$\rho \left( \frac{\partial u^i}{\partial t} + u^j \frac{\partial u^i}{\partial x^j} \right) = \frac{-\partial p}{\partial x^i} + (\lambda + \mu) \frac{\partial}{\partial x^i} \text{div} u + \mu \Delta u^i$$

If  $\mu = 0$ ,  $\text{div} u = 0$  and  $\frac{\partial u^i}{\partial t} = 0$  (steady state),

$$\rho_0 u^j \frac{\partial u^i}{\partial x^j} + \frac{\partial p}{\partial x^i} = 0$$

For  $j = i$ ,  $\frac{\partial}{\partial x^i} (\frac{1}{2} \rho_0 (u^i)^2 + p) + \sum_{j \neq i} \rho_0 u^j \frac{\partial u^i}{\partial x^j} = 0$

Part 11. Electromagnetism

cf. 1-2 Exhaust velocity and Specific Impulse of Jahn (2006) [27]  
Jahn (2006) [27] defines specific impulse  $I_s$  as the ratio of thrust to the rate of use of propellant by sea level weight:

(231) 
$$I_s := \frac{T}{\dot{m} g_0}$$

**Problem 1-4.** It should be pointed out that the physics is exactly the same as before as propelling mass out, *before* separation.  
What’s different is the final mass after all planned for propellant mass is expelled is different.  
To review,

Let  $M = M(t)$  the mass of the entire system at any time, only excluding propellant mass that’s *already* expelled out.

$M\dot u = -\dot M u_{\text{exh}} + F_g$ . For the case of  $F_g = -Mg$ ,

$$\begin{aligned}\dot u &= \frac{-\dot M}{M} u_{\text{exh}} + \frac{F_g}{M} = \frac{-\dot M}{M} u_{\text{exh}} + -g \\ \xrightarrow{\int dt} \Delta u &\equiv u(t_f) - u(0) = -\ln\left(\frac{M(t_f)}{M(0)}\right) u_{\text{exh}} - g t_f \\ \implies \frac{\Delta u}{u_{\text{exh}}} + \frac{g \Delta t}{u_{\text{exh}}} &= \ln\left(\frac{M(0)}{M(t_f)}\right) \text{ or } \exp\left(\frac{\Delta u}{u_{\text{exh}}}\right) \exp\left(\frac{g \Delta t}{u_{\text{exh}}}\right) = \frac{m_0}{M(t_f)}\end{aligned}$$

where  $t_f \equiv \Delta t$  and  $m_0 \equiv M(0)$  (Jahn’s notation in the former, my notation for the latter).  
Now the mass of the power supply,  $m_p$  was given by Jahn and reasoned to monotonically depend on the power  $P$ , which is reasonable. Also, we have  $M(0)$ , and  $M(t_f) \equiv M(\Delta t)$  values:

$$\begin{aligned}m_p &= \alpha P = \alpha \frac{T u_{\text{exh}}}{2\eta} = \frac{\alpha \dot m u_{\text{exh}}^2}{2\eta} \\ M(\Delta t) &= m_u + m_p = m_u + \frac{\alpha}{2\eta} \frac{M_p}{\Delta t} u_{\text{exh}}^2 \\ m_0 \equiv M(0) &= M_p + m_u + m_p = M_p + m_u + \frac{\alpha \dot m u_{\text{exh}}^2}{2\eta} = \\ &= M_p + m_u + \alpha \left(\frac{M_p}{\Delta t}\right) \frac{u_{\text{exh}}^2}{2\eta} \\ \implies M_p &= \frac{m_0 - m_u}{1 + \frac{\alpha u_{\text{exh}}^2}{2\eta \Delta t}}\end{aligned}$$

for  $\dot m \Delta t = M_p$  (*assume constant* mass flow out, and by  $\Delta t$  time,  $M_p$  mass propellant is used up).  
Plugging all this in:

$$\exp\left(\frac{\Delta u}{u_{\text{exh}}}\right) \exp\left(\frac{g \Delta t}{u_{\text{exh}}}\right) = \frac{m_0}{m_u + \frac{\alpha}{2\eta \Delta t} \left(\frac{m_0 - m_u}{1 + \frac{\alpha u_{\text{exh}}^2}{2\eta \Delta t}}\right) u_{\text{exh}}^2} = \frac{1 + \frac{\alpha u_{\text{exh}}^2}{2\eta \Delta t}}{\frac{m_u}{m_0} + \frac{\alpha u_{\text{exh}}^2}{2\eta \Delta t}}$$

Let  $g = 0$ . Then the desired result is obtained, for ratio of mass payload to initial mass, including a power supply:

(232) 
$$\frac{m_u}{m_0} = e^{-\Delta u / u_{\text{exh}}} + \frac{\alpha u_{\text{exh}}^2}{2\eta \Delta t} (e^{-\Delta u / u_{\text{exh}}} - 1)$$

38.1. **Charge conservation.** Recall that the divergence div:

$$\begin{aligned}\text{div} : \mathfrak{X}(M) &\rightarrow C^\infty(M) \\ \text{div} \mathbf{j} &= \frac{1}{\sqrt{g}} \frac{\partial(j^k \sqrt{g})}{\partial x^k}\end{aligned}$$

Usually, charge conservation is written in differential vector form as

$$\frac{\partial \rho}{\partial t} + \text{div} \mathbf{j} = 0$$

Does this hold for the differential form version of electromagnetism?  
Consider  $j \in \Omega^1(M)$ . Now

$$j = j_\mu dx^\mu = \rho dt + j_i dx^i \in \Omega^1(M) \overset{\#}{\mapsto} g^{\mu\nu} j_\nu \frac{\partial}{\partial x^\mu} \equiv j^\mu \frac{\partial}{\partial x^\mu} \in \mathfrak{X}(M)$$

Then

$$\begin{aligned} *j &= \frac{\sqrt{g}}{3!} \epsilon_{\mu\nu\rho\sigma} g^{\mu\mu_1} j_{\mu_1} dx^\nu \wedge dx^\rho \wedge dx^\sigma \\ d *j &= \frac{1}{3!} \epsilon_{\mu\nu\rho\sigma} \frac{\partial}{\partial x^{\mu_2}} (\sqrt{g} g^{\mu\mu_1} j_{\mu_1}) dx^{\mu_2} \wedge dx^\nu \wedge dx^\rho \wedge dx^\sigma = \frac{1}{3!} \epsilon_{\mu\nu\rho\sigma} \frac{\partial}{\partial x^{\mu_2}} (\sqrt{g} j^\mu) dx^{\mu_2} \wedge dx^\nu \wedge dx^\rho \wedge dx^\sigma = \\ &= \frac{1}{3!} \{ \epsilon_{ijk} \frac{\partial(\sqrt{g} j^0)}{\partial t} dt \wedge dx^i \wedge dx^j \wedge dx^k + \epsilon_{ijk} \frac{\partial}{\partial x^l} (\sqrt{g} j^i) dx^l \wedge dx^j \wedge dx^k + \dots \} = 0 \\ \implies \frac{1}{\sqrt{g}} \frac{\partial(\sqrt{g} j^0)}{\partial t} + \text{div} \mathbf{j} &= \frac{1}{\sqrt{g}} \frac{\partial(\sqrt{g} g^{0\mu} j_\mu)}{\partial t} + \frac{1}{\sqrt{g}} \frac{\partial(\sqrt{g} g^{k\mu} j_\mu)}{\partial x^k} = 0\end{aligned}$$

Suppose current is a **convective current**  $\mathbf{j} = \rho \mathbf{u}$ , where  $\mathbf{u} \in \mathfrak{X}(M)$  ( $\mathbf{u} = \mathbf{u}(t, x)$ ,  $(t, x) \in M$ ).

Charge conservation always holds.  
my version of Maxwell’s equations:  
if  $\nabla \cdot \mathbf{B} = 0$ ,  
then  $\nabla \times \mathbf{E} = \frac{-1}{c} \left(\frac{\partial \mathbf{B}}{\partial t}\right)$   
if  $\nabla \cdot \mathbf{E} = 4\pi \rho_{\text{total}}$   
then  $\nabla \times \mathbf{B} = \frac{1}{c} \left(\frac{\partial \mathbf{E}}{\partial t} + 4\pi \frac{\partial \mathbf{P}}{\partial t} + 4\pi \mathbf{J}_{\text{free}} + 4\pi c \nabla \times \mathbf{M}\right)$   
cf. Section 1.4 Coulomb’s law of Purcell (1984) [28].  
For  $\mathbf{F}_2$  = force on charge 2,

$$\mathbf{F}_2 = \frac{k q_1 q_2 \mathbf{r}_{21}}{r_{21}^2} = \frac{k q_1 q_2 (\mathbf{r}_2 - \mathbf{r}_1)}{|\mathbf{r}_2 - \mathbf{r}_1|^3}$$

I am curious to know if the tangent-cotangent (“musical”) isomorphism holds:

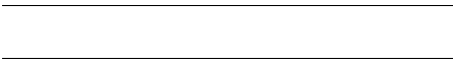
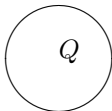
$$F_2^i = \frac{k q_1 q_2 (\mathbf{r}_2 - \mathbf{r}_1)^i}{|\mathbf{r}_2 - \mathbf{r}_1|^3} \in \mathfrak{X}(M) \overset{\flat}{\mapsto} (F_2)_j = \frac{k q_1 q_2 (\mathbf{r}_2 - \mathbf{r}_1)^i g_{ij}}{|\mathbf{r}_2 - \mathbf{r}_1|^3} \in \Omega^1(M)$$

$k = 1$  (cgs)  
 $k = \frac{1}{4\pi\epsilon_0}$ ;  $\epsilon_0 = 8.85 \times 10^{-12}$  ;  $k = 8.988 \times 10^9$  (SI)  
1 coulomb =  $2.998 \times 10^9$  esu  $\cong 3 \times 10^9$  esu  
 $\epsilon_0 = 4.8 \times 10^{-16}$  esu =  $1.6 \times 10^{-16}$  C  
 $10^5$  dynes = 1 Newtons.  
So the above takes care of Ch. 2 Electromagnetic Theory, Section 2-1 Electric Charges and Electrostatic Fields of Jahn (2006) [27] up to Eqn. (2-1).  
Section 3.5 Capacitance and Capacitors of Purcell (1984) [28].  
Consider an isolated conductor carrying charge  $Q$ .

$$Q = C \phi_0$$

cgs.  $[Q] = \text{esu}$ ,  $[C] = \text{cm}$ ,  $[\phi_0] = \text{statvolts} = \frac{\text{ergs}}{\text{esu}} = \frac{\text{dyne}\cdot\text{cm}}{\text{esu}}$ .  
SI.  $[Q] = \text{coulombs}$ ,  $[C] = \text{farads}$ ,  $[\phi] = \text{volts}$ .





cgs:  $C = a$   $[a] = m$   $C = \frac{A}{4\pi s}$   $[A] = \text{cm}^2$ ,  $[s] = \text{cm}$   
SI:  $C = 4\pi\epsilon_0 a$   $[a] = m$   $C = \frac{\epsilon_0 A}{s}$   $[A] = \text{m}^2$ ,  $[s] = \text{m}$   
 $\epsilon_0 = 8.854 \times 10^{-12}$  farad/meter

1 cm = 1.11 × 10<sup>−12</sup> farad

Fitzpatrick (2014) [29]

Part 12. Physical Kinetics

Pitaevskii and Lifshitz (1981) [30]  
Sonnendrücker (2017) [33].

Part 13. Guidance, Navigation, and Control (GNC)

39. SOFTWARE FOR MODELING AND EXECUTION

39.1. Matlab/Simulink and alternatives.

39.1.1. *Open source alternatives.* GNURadio for digital signal processing (cf. Nick C.)  
In image processing, Matlab replaced with Python (openCV, numpy, PIL, other great modules) (cf. Katharine I.)  
Modelica - has fairly good documentation with examples, has standard library that may or may not be surface level of domain of interest. Has FMI support, allowing saving models/simulations in standardized plain text data. (cf. Sean M.)  
Use of Matlab is very team oriented (meaning very specific to the team’s needs?) For example for Power Electronics, LTSpice is used, as well for PCB simulations. Simulink is good for model building and doing complex systems. (cf. Amir K.)

Part 14. Control, Control Theory

40. SIGNALS AND SYSTEMS

cf. pp. 64, Oppenheim and Willsky (1997) [44]

**Problem 1.35.** Consider periodic discrete-time exponential time signal:

$$x[n] = e^{im(2\pi/N)n}$$

**Solution 1.35.**  
For  $\exp\left(i\frac{2\pi m}{N}n\right)$ ,  $n \in \mathbb{Z}$ , consider  $\frac{2\pi m}{N} = \frac{2\pi}{T_0}$  or  $T_0 = N/m$ .  
Consider  $\gcd(N, m) = d$ , then  $\exists r, s \in \mathbb{Z}$ , s.t.  $N = dr$ ,  $m = ds$ , and  $\gcd(r, s) = 1$ .

$$\begin{aligned}\frac{N}{\gcd(m, N)} &= \frac{N}{d} = \frac{N}{N/r} = r \\ \frac{N}{d} &= \frac{N}{m/s} = \frac{sN}{m} \\ r &= \frac{sN}{m} \text{ or } \frac{r}{s} = \frac{N}{m}\end{aligned}$$

But  $\gcd(r, s) = 1$ , whereas  $\gcd(m, N) = d$ , so  $\frac{N}{m} = \frac{N/d}{m/d}$

**Problem 1.36.**

- a Let  $x(t)$  cont.-time complex exponential signal  $x(t) = e^{i\omega_0 t}$   
(1)  
c Assuming  $\frac{T}{T_0} = \frac{p}{q}$  (P1.36-1, pp. 64, Oppenheim and Willsky (1997) [44]), determine precisely how many periods of  $x(t)$  are needed to obtain samples that form single period of  $x[n]$ .

**Solution 1.36.**

(a)

$$\text{for } \exp(i\omega_0 nT), \quad \omega_0(n+k)T = \omega_n nT + 2\pi l \implies \omega_0 kT = 2\pi l$$
$$\frac{T}{(2\pi/\omega_0)} = \frac{l}{k} \text{ where } l, k \in \mathbb{Z}$$

We had let  $T_0 = 2\pi/\omega_0$

Jacobs (1994) [45]

41. INTRODUCTION TO CONTROL SYSTEMS

cf. Ch. 1 of Ogata (2009) [47]  
cf. Sec. 1.1 ”Introduction” of Ogata (2009) [47]  
**Plants** - plant maybe piece of equipment, the purpose of which is to perform a particular operation. Any physical object to be controlled (e.g. mechanical device, heating furnance, chemical reactor, spacecraft) is a plant.  
cf. 1-3 ”Closed loop Control versus Open-loop control” of Ogata (2009) [47].  
**Feedback Control Systems** system that maintains prescribed relationship between output and reference input by comparing them and using the difference as means of control is a *feedback control system*.  
**Closed-Loop Control Systems** In practice, terms feedback control and closed-loop control are used interchangeably.In a closed-loop control system, the actuating error signal, which is the difference between input signal and feedback signal (which may be the output signal itself or function of the output signal and its derivatives and/or integrals), is fed to controller so to reduce error and bring output of system to a desired value.  
**Open-Loop control systems.** Those systems in which output has no effect on control action are called *open-loop control systems*.

42. MATHEMATICAL MODELING OF CONTROL SYSTEMS, BLOCK DIAGRAMS

42.1. **Transfer Function.** cf. Ch. 2 ”Mathematical Modeling of Control Systems”, Ogata (2009) [47]  
Consider lineaer time-invariant sytem:

(233)

$$\sum_{l=0}^{n-1} a_l \frac{d^{n-l}y}{dt^{n-l}} + a_n y = \sum_{k=0}^{m-1} b_k \frac{d^{m-k}x}{dt^{m-k}} + b_m x \quad (n \geq m)$$

where  $y$  is the output of the system, and  $x$  is input.

The **transfer function** of this system is the ratio of the Laplace transformed output to the Laplace transformed input when all initial conditions are zero, or

(234)

$$G(s) = \frac{\mathcal{L}[y]}{\mathcal{L}[x]} \Big|_{x(0)=y(0)=0} = \frac{Y(s)}{X(s)} = \frac{\sum_{k=0}^{m-1} b_k s^{m-k} + b_m}{\sum_{l=0}^{n-1} a_l s^{n-l} + a_n}$$

**Comments on Transfer Function.** The applicability of the concept of the transfer function is limited to linear, time-invariant, differential equation systems.

- (1)
- (2) The transfer function is a property of a system itself, independent of the magnitude and nature of the input or driving function.
- (3) The transfer function includes the units necessary to relate input to output; however, it doesn't provide any information concerning the physical structure of the system. (The transfer functions of many physically different systems can be identical.)
- (4)
- (5)

TODO: Proof with Fubini's, cf. [Math stackexchange, "Proof of convolution theorem for Laplace transform" https://www.sciencedirect.com/topics/computer-science/convolution-theorem](https://www.sciencedirect.com/topics/computer-science/convolution-theorem)

$$E(s) = R(s) - B(s) \text{ (at summing point)}$$

Feedback element (accomplishes conversion, modify output before it's compared with input), whose transfer function is  $H(s)$ .

(235)

$$H(s) = \frac{B(s)}{C(s)} \text{ or } B(s) = H(s)C(s)$$
$$G(s) = \frac{C(s)}{E(s)}$$

**42.2. Automatic Control Systems, Block Diagrams.** cf. Sec. 2-3 "Automatic Control Systems", pp. 17, Ogata (2009) [\[47\]](#)

**42.2.1. Open-loop Transfer Function and Feedforward Transfer Function.** Ratio of feedback signal  $B(s)$  to actuating error signal  $E(s)$  is **open-loop transfer function**

(236)

$$\text{Open-loop transfer function} = \frac{B(s)}{E(s)} = G(s)H(s)$$

ratio of output  $C(s)$  to actuating error signal  $E(s)$  is feedforward transfer function so

(237)

$$\text{Feedforward transfer function} = \frac{C(s)}{E(s)} = G(s)$$

**42.2.2. Closed Loop Transfer Function.** Now

$$C(s) = G(s)E(s)$$
$$E(s) = R(s) - B(s) = R(s) - H(s)C(s)$$
$$\implies C(s) = G(s)(R(s) - H(s)C(s)) \text{ or}$$

(238)

$$\frac{C(s)}{R(s)} = \frac{G(s)}{1 + G(s)H(s)}$$

Transfer function relating  $C(s)$  to  $R(s)$  is called *closed-loop transfer function*. It relates closed-loop system dynamics to dynamics of feedforward elements and feedback elements.

**42.2.3. Closed Loop System Subjected to a Disturbance.** Recall

$$E(s) = R(s) - B(s) = R(s) - H(s)C(s)$$

When 2 inputs (reference input and disturbance) are present in a linear time-invariant system, each input can be treated independently of the other.

In examining effect of disturbance  $D(s)$ , assume reference input 0

$$G_1(s) = \frac{Y_1(s)}{-B(s)}$$
$$G_2(s) = \frac{C_D(s)}{D(s) + Y_1(s)} = \frac{C_D(s)}{D(s) + -G_1(s)B(s)} = \frac{C_D(s)}{D(s) - G_1(s)(H(s)C_D(s))}$$
$$G_2(s)D(s) - G_2(s)G_1(s)H(s)C_D(s) = C_D(s)$$

so

(239)

$$\frac{C_D(s)}{D(s)} = \frac{G_2(s)}{1 + G_1(s)G_2(s)H(s)}$$

In considering response to reference input, assume disturbance 0. Then response  $C_R(s)$  to reference input  $R(s)$  is given by

$$E(s) = R(s) - B(s) = R(s) - H(s)C(s)$$

$$G_1(s) = \frac{Y(s)}{E(s)}$$
$$G_2(s) = \frac{C_R(s)}{Y(s)} \implies G_1(s)G_2(s) = \frac{C_R(s)}{E(s)} = \frac{C_R(s)}{R(s) - H(s)C_R(s)}$$

thus

(240)

$$\frac{C_R(s)}{R(s)} = \frac{G_1(s)G_2(s)}{1 + G_1(s)G_2(s)H(s)}$$

Response to simultaneous application of reference input and disturbance obtained by adding:

(241)

$$C(s) = C_R(s) + C_D(s)$$
$$\frac{G_2(s)}{1 + G_1(s)G_2(s)H(s)} [G_1(s)R(s) + D(s)]$$

Consider when  $|G_1(s)H(s)| \gg 1$  and  $|G_1(s)G_2(s)H(s)| \gg 1$ .

**42.3. Modeling in State Space.** cf. pp. 29 "Modeling in State Space", Sec. 2-4, Ogata (2009) [\[47\]](#)  
Define

$$\mathbf{x}(t) = (\mathbf{x}(t))_i = x_i(t), \quad 1 \leq i \leq n$$
$$\mathbf{y}(t) = (\mathbf{y}(t))_i = y_i(t), \quad 1 \leq i \leq m$$
$$\mathbf{u}(t) = (\mathbf{u}(t))_i = u_i(t), \quad 1 \leq i \leq r$$
$$\mathbf{f}(\mathbf{x}, \mathbf{u}, t) = (\mathbf{f}(\mathbf{x}, \mathbf{u}, t))_i = f_i(x_1, x_2, \dots x_n; u_1, u_2, \dots u_r, t), \quad 1 \leq i \leq n$$
$$\mathbf{g}(\mathbf{x}, \mathbf{u}, t) = (\mathbf{g}(\mathbf{x}, \mathbf{u}, t))_i = g_i(x_1, x_2, \dots x_n; u_1, u_2, \dots u_r, t), \quad 1 \leq i \leq m$$

Then

(243)

$$\dot{\mathbf{x}}(t) = \mathbf{f}(\mathbf{x}, \mathbf{u}, t)$$
$$\mathbf{y}(t) = \mathbf{g}(\mathbf{x}, \mathbf{u}, t)$$

If Eq. 243 is linearized about the operating state, then we have the following linearized state equation and output equation:

$$(244) \quad \begin{aligned} \dot{\mathbf{x}}(t) &= \mathbf{A}(t)\mathbf{x}(t) + \mathbf{B}(t)\mathbf{u}(t) \\ \mathbf{y}(t) &= \mathbf{C}(t)\mathbf{x}(t) + \mathbf{D}(t)\mathbf{u}(t) \end{aligned}$$

where  $\mathbf{A}(t) \equiv$  state matrix,  $\mathbf{B}(t) \equiv$  input matrix,  $\mathbf{C}(t) \equiv$  output matrix,  $\mathbf{D}(t) \equiv$  direct transmission matrix.

42.3.1. *Correlation Between Transfer Functions and State-Space Equations.* cf. "Correlation Between Transfer Functions and State-Space Equations", pp. 33, Sec. 2-4 "Modeling in State Space", Ogata (2009) [47].

Show how to derive transfer function of single-input, single-output system from state-space equations.

Consider system with transfer function

$$\frac{Y(s)}{U(s)} = G(s)$$

System in state space representation:

$$\begin{aligned} \dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}u \\ y &= \mathbf{C}\mathbf{x} + Du \end{aligned}$$

where  $\mathbf{x}$  is state vector,  $u$  is input,  $y$  is output.

Apply Lapalce Transforms:

$$(245) \quad \begin{aligned} s\mathbf{X}(s) - \mathbf{X}(0) &= \mathbf{A}\mathbf{X}(s) + \mathbf{B}U(s) \\ Y(s) &= \mathbf{C}\mathbf{X}(s) + DU(s) \end{aligned}$$

Let  $\mathbf{X}(0) = 0$ . Then  $s\mathbf{X}(s) - \mathbf{A}\mathbf{X}(s) + \mathbf{B}U(s)$  or  $(s\mathbf{1} - \mathbf{A})\mathbf{X}(s) = \mathbf{B}U(s)$ ; premultiply by  $(s\mathbf{1} - \mathbf{A})^{-1}$  to get

$$\begin{aligned} \mathbf{X}(s) &= (s\mathbf{1} - \mathbf{A})^{-1}\mathbf{B}U(s) \\ \implies Y(s) &= [\mathbf{C}(s\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + D]U(s) \end{aligned}$$

so

$$(246) \quad G(s) = \mathbf{C}(s\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + D$$

Since RHS of  $G(s)$  in Eq. 246 involces  $(s\mathbf{1} - \mathbf{A})^{-1}$ ,  $G(s)$  can be written as

$$G(s) = \frac{Q(s)}{|s\mathbf{1} - \mathbf{A}|}$$

where  $Q(s)$  is a polynomial in  $s$ .

Notice  $|s\mathbf{1} - \mathbf{A}| =$  characteristic polynomial of  $G(s)$ , i.e. eigenvalues of  $\mathbf{A} =$  poles of  $G(s)$ .

42.3.2. *Transfer Matrix.* Consider multiple input, multiple-output system. Assume  $r$  inputs  $u_1, u_2, \dots, u_r$ ,  $m$  outputs  $y_1, y_2, \dots, y_m$ . Define  $\mathbf{y} = y_i$ ,  $1 \leq i \leq m$ ,  $\mathbf{u} = u_i$ ,  $1 \leq i \leq r$ .

Transfer matrix  $\mathbf{G}(s)$  relates output  $\mathbf{Y}(s)$  to input  $\mathbf{U}(s)$  or  $\mathbf{Y}(s) = \mathbf{G}(s)\mathbf{U}(s)$ , where

$$(247) \quad \mathbf{G}(s) = \mathbf{C}(s\mathbf{1} - \mathbf{A})^{-1}\mathbf{B} + \mathbf{D}$$

#### 42.4. State-Space Representation of Scalar Differential Equation Systems.

42.4.1. *State-Space Representation of  $n$ th-order systems of Linear Differential Equations in which the Forcing Function Does Not Involve Derivative Terms.* Consider

$$(248) \quad \frac{d^n y}{dt^n} + \sum_{l=1}^{n-1} a_l \frac{d^{n-l} y}{dt^{n-l}} + a_n y = u$$

Noting that knowledge of  $y(0), \dot{y}(0), \dots, \frac{d^{n-1}}{dt^{n-1}}y(0)$ , and input  $u(t)$  for  $t \geq 0$  completely determines system, take  $\{y(t), \dot{y}(t), \dots, \frac{d^{n-1}}{dt^{n-1}}y(t)\}$  as set of  $n$  state variables.

Define

$$(249) \quad \begin{aligned} x_1 &= y \\ x_i &= \frac{d^{i-1} y}{dt^{i-1}} \quad 2 \leq i \leq n \end{aligned}$$

Then

$$(250) \quad \begin{aligned} \dot{x}_{i-1} &= x_i, \quad 2 \leq i \leq n \\ \dot{x}_n &= -\sum_{l=1}^n a_{n+1-l} x_l + u \end{aligned}$$

or

$$(251) \quad \begin{aligned} \dot{\mathbf{x}} &= \mathbf{A}\mathbf{x} + \mathbf{B}u \text{ where} \\ \mathbf{x} &= (\mathbf{x})_i = x_i \quad 1 \leq i \leq n \text{ where} \\ \mathbf{A} &= (\mathbf{A})_{ij} = \begin{cases} \delta_{i,j-1} & \text{if } i \neq n \\ -a_{n-j} & \text{if } i = n \end{cases} \quad \mathbf{B} = \delta_{i,n} \end{aligned}$$

$$(252) \quad y = [1 \quad 0 \quad \dots \quad 0] \mathbf{x} = \mathbf{C}\mathbf{x}$$

Note that state-space representation for transfer function system:

$$(253) \quad \frac{Y(s)}{U(s)} = \frac{1}{s^n + \sum_{l=1}^n a_l s^{n-l}}$$

42.4.2. *State-Space Representation of  $n$ th-order systems of Linear Differential Equations in which the Forcing Function Involves Derivative Terms.* cf. pp. 37, Sec. 2-5 "State-Space Representation of Scalar Differential Equation Systems", Ogata (2009) [47]

Consider system involving derivatives of forcing function  $u$ :

$$(254) \quad \frac{d^n y}{dt^n} + \sum_{l=1}^{n-1} a_l \frac{d^{n-l} y}{dt^{n-l}} + a_n y = \sum_{l=0}^{n-1} b_l \frac{d^{n-l} u}{dt^{n-l}} + b_n u$$

Main problem in defining state variables for this case is derivative terms of input  $u$ . State variables must be such that they'll eliminate derivatives of  $u$  in state equation.

Try

$$(255) \quad \begin{aligned} x_1 &= y - \beta_0 u \\ x_2 &= \dot{y} - \beta_0 \dot{u} - \beta_1 u = \dot{x}_1 - \beta_1 u \\ x_i &= \frac{d^{i-1} y}{dt^{i-1}} - \sum_{j=0}^{i-2} \beta_j \frac{d^{i-1-j} u}{dt^{i-1-j}} - \beta_{i-1} u = \dot{x}_{i-1} - \beta_{i-1} u \quad 2 \leq i \leq n \end{aligned}$$

cf. (2-34) of Ogata (2009) [47], where  $\beta_j$ ,  $0 \leq j \leq n-1$  determined by

(256)

$$\begin{aligned}\beta_0 &= b_0 \\ \beta_1 &= b_1 - a_1\beta_0 \\ \beta_i &= b_i - \sum_{l=1}^i a_l\beta_{n-1-l} \quad 1 \leq i \leq n-1\end{aligned}$$

cf. (2-35) of Ogata (2009) [47], with this choice of state variables, existence and uniqueness of solution of state equation is guaranteed. (EY ???) With the present choice of state variables, we obtain (derivation in Problem A-2-6):

(257)

$$\begin{aligned}\dot{x}_i &= x_{i+1} + \beta_i u \quad 1 \leq i \leq n-1 \\ \dot{x}_n &= -\sum_{l=1}^n a_{n+1-l}x_l + \beta_n u\end{aligned}$$

cf. (2-36) of Ogata (2009) [47], where

$$\beta_n = b_n - \sum_{l=1}^n a_l\beta_{n-l}$$

**Problem A-2-6.** From the definition of state variables, Eq. 255,

$$x_{i+1} = \dot{x}_i - \beta_i u \text{ or } \dot{x}_i = x_{i+1} + \beta_i u \text{ since}$$

$$\begin{aligned}\dot{x}_i &= \frac{d^i y}{dt^i} - \sum_{j=0}^{i-2} \beta_j \frac{d^{i-j} u}{dt^{i-j}} - \beta_{i-1} \dot{u} = \frac{d^i y}{dt^i} - \sum_{j=0}^{i-2} \beta_j \frac{d^{i-j} u}{dt^{i-j}} - \beta_{i-1} \dot{u} - \beta_i u + \beta_i u \\ &= \frac{d^i y}{dt^i} - \sum_{j=0}^{i-1} \beta_j \frac{d^{i-j} u}{dt^{i-j}} - \beta_i u + \beta_i u = x_{i+1} + \beta_i u\end{aligned}$$

Note that from Eq. 254

$$\frac{d^n y}{dt^n} = -\sum_{l=1}^{n-1} a_l \frac{d^{n-l} y}{dt^{n-l}} - a_n y + \sum_{l=0}^{n-1} b_l \frac{d^{n-l} u}{dt^{n-l}} + b_n u$$

Now since

$$x_n = \frac{d^{n-1} y}{dt^{n-1}} - \sum_{j=0}^{n-2} \beta_j \frac{d^{n-1-j} u}{dt^{n-1-j}} - \beta_{n-1} u$$

Then

$$\dot{x}_n = -\sum_{l=1}^{n-1} a_l \frac{d^{n-l} y}{dt^{n-l}} - a_n y + \sum_{l=0}^{n-1} b_l \frac{d^{n-l} u}{dt^{n-l}} + b_n u - \sum_{j=0}^{n-2} \beta_j \frac{d^{n-j} u}{dt^{n-j}} - \beta_{n-1} \dot{u}$$

Substitute in for  $\frac{d^{n-l} y}{dt^{n-l}}$  from Eq. 255:

$$\dot{x}_n = -\sum_{l=1}^{n-1} a_l \left[ x_{n-l+1} + \sum_{j=0}^{n-l-1} \beta_j \frac{d^{n-l-j} u}{dt^{n-l-j}} + \beta_{n-l} u \right] - a_n y + \sum_{l=0}^{n-1} b_l \frac{d^{n-l} u}{dt^{n-l}} + b_n u - \sum_{j=0}^{n-1} \beta_j \frac{d^{n-j} u}{dt^{n-j}} =$$

Consider

$$\sum_{l=0}^{n-1} (b_l - \beta_l) \frac{d^{n-l} u}{dt^{n-l}} = \sum_{l=0}^{n-1} \left( \sum_{j=1}^l a_j \beta_{n-1-j} \right) \frac{d^{n-l} u}{dt^{n-l}}$$

Consider the double summation in the terms

$$\begin{aligned}-\sum_{l=1}^{n-1} a_l \sum_{j=0}^{n-l-1} \beta_j \frac{d^{n-l-j} u}{dt^{n-l-j}} \text{ and} \\ \sum_{l=0}^{n-1} \left( \sum_{j=1}^l a_j \beta_{n-1-j} \right) \frac{d^{n-l} u}{dt^{n-l}}\end{aligned}$$

For each, plot the index  $l$  on the  $x$ -axis (independent), and  $j$  on the  $y$  axis (dependent). For each double summation, notice the valid  $(j, l)$  indices form a triangle and each double summation have the same number of points on this  $l-j$  grid.

First, do this, adding one to the index of  $j$  for this double sum:

$$\sum_{l=1}^{n-1} \sum_{j=0}^{n-l-1} = \sum_{l=1}^{n-1} \sum_{j=1}^{n-l}$$

Consider making the following substitutions:

$$\sum_{l=1}^{n-1} \sum_{j=1}^{n-l} \xrightarrow{m=n-1-l} \sum_{m=0}^{n-2} \sum_{j=1}^{m+1} = \sum_{m=1}^{n-1} \sum_{j=1}^m$$

So then

$$\begin{aligned}\dot{x}_n &= -\sum_{l=1}^{n-1} a_l x_{n-l+1} - \sum_{l=1}^{n-1} a_l \beta_{n-l} u - a_n y + b_n u = -\sum_{l=1}^{n-1} a_l x_{n-l+1} + \left( b_n - \sum_{l=1}^{n-1} a_l \beta_{n-l} \right) u - a_n (x_1 + \beta_0 u) \\ &= -\sum_{l=1}^n a_l x_{n-l+1} + \beta_n u\end{aligned}$$

where  $x_1 = y - \beta_0 u$  and then  $\beta_n = b_n - \sum_{l=1}^n a_l \beta_{n-l}$  was used.

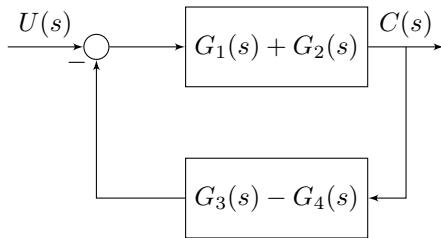
**Problem B-2-1.** cf. pp. 60, Ch. 2 ”Mathematical Modeling of Control Systems”, Ogata (2009) [47]

$$\begin{aligned}G_1 &= \frac{C_1}{E} \\ G_2 &= \frac{C_2}{E} \\ G_3 &= \frac{Y_3}{C} \\ G_4 &= \frac{Y_4}{C}\end{aligned}$$

Then

$$\begin{aligned}C_1 + C_2 &= C = (G_1 + G_2)E = (G_1 + G_2)(R - C(G_3 - G_4)) \\ E &= R - (Y_3 - Y_4)\end{aligned}$$

$$\frac{C}{R} = \frac{G_1 + G_2}{1 + (G_1 + G_2)(G_3 - G_4)}$$



**Problem B-2-7 of Ogata (2009) [47].** pp. 61, Ch. 2 "Mathematical Modeling of Control Systems"  
If

$$G_c(s) = \frac{Y_c(s)}{E(s)}$$

$$G_p(s) = \frac{Y_p(s)}{Y_c(s)}$$

$$\text{so } G_c(s)G_p(s) = \frac{Y_p(s)}{E(s)}$$

Suppose  $R(s) = 0$ . Let the response be  $C_D$ .

$$G_c(s)G_p(s) = \frac{Y_p(s)}{-C_D(s)}$$

$$Y_p(s) + D(s) = C_D(s)$$

$$\implies -G_c(s)G_p(s)C_D(s) = C_D(s) - D(s)$$

so

$$\boxed{\frac{C_D(s)}{D(s)} = \frac{1}{1 + G_c(s)G_p(s)}}$$

Now let  $D(s) = 0$ .

$$G_c(s)G_p(s) = \frac{C_R(s)}{R(s) - C_R(s)} \text{ or } G_c(s)G_p(s)(R(s) - C_R(s)) = C_R(s)$$

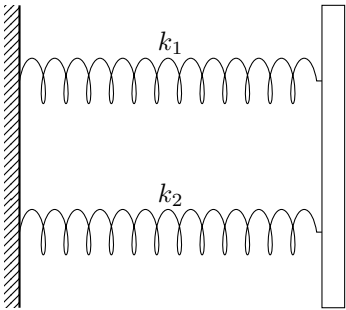
$$\frac{C_R(s)}{R(s)} = \frac{G_c(s)G_p(s)}{1 + G_c(s)G_p(s)}$$

$$\implies \boxed{C(s) = C_R(s) + C_D(s) = \frac{R(s)G_c(s)G_p(s) + D(s)}{1 + G_c(s)G_p(s)}}$$

#### 43. MATHEMATICAL MODELING OF MECHANICAL SYSTEMS

Example 3-1 of Ogata (2009) [47] **(Springs in Parallel and Series)**

*Parallel:*



$$F = -k_1x - k_2x = -(k_1 + k_2)x = -k_{\text{eq}}x$$

For springs in series, the *force in each spring is the same* (EY??).

$$F = -k_1x_1 \quad F = -k_2(x_2 - x_1) \text{ so } F = -k_2x_2 + k_2\left(\frac{-F}{k_1}\right) \text{ or } F(1 + \frac{k_2}{k_1}) = -k_2x_2$$

$$\implies \frac{F}{-x_2} = k_{\text{eq}} = \frac{k_1k_2}{k_1 + k_2} = \frac{1}{\frac{1}{k_1} + \frac{1}{k_2}}$$

Example 3-2 of Ogata (2009) [47] **(Dashpots in Parallel and Series)**

*Parallel:* Force  $f$  due to dampers:

$$f = -b_1(\dot{x}_2 - \dot{x}_1) - b_2(\dot{x}_2 - \dot{x}_1) = -(b_1 + b_2)(\dot{x}_2 - \dot{x}_1)$$

or

$$f = -b_{\text{eq}}(\dot{x}_2 - \dot{x}_1) \text{ where } b_{\text{eq}} = b_1 + b_2$$

**Example 3-3.** Consider spring-mass-dashpot system mounted on massless cart. Assume cart standing still for  $t < 0$ , and spring-mass-dashpot system on cart also standing still for  $t < 0$ .  
 $y(t) \equiv$  displacement of mass is output (relative to ground),  $m \equiv$  mass,  $b \equiv$  viscous-friction coefficient,  $k \equiv$  spring constant

$$(258) \quad ma = \sum F \text{ so } m\ddot{y} = -b(\dot{y} - \dot{u}) - k(y - u) \text{ or } m\ddot{y} + b\dot{y} + ky = b\dot{u} + ku$$

Doing the Laplace transform:

$$(ms^2 + bs + k)Y(s) = (bs + k)U(s)$$

Thus we obtain the *transfer function*

$$(259) \quad G(s) = \frac{Y(s)}{U(s)} = \frac{bs + k}{ms^2 + bs + k}$$

Such a transfer function representation of a mathematical model is used very frequently in control engineering.  
Next, obtain state-space model of this system.  
Rewrite dynamics in Eq. 258 as

$$(260) \quad \ddot{y} + \frac{b}{m}\dot{y} + \frac{k}{m}y = \frac{b}{m}\dot{u} + \frac{k}{m}u = \ddot{y} + a_1\dot{y} + a_2y = b_1\dot{u} + b_2u$$

Now  $\beta_0 = b_0 = 0$ ,  $\beta_1 = b_1 - a_1\beta_0 = \frac{b}{m}$ ,  $\beta_2 = \frac{k}{m} - \left(\frac{b}{m}\right)^2$ .  
From Eq. 256,

$$\beta_0 = b_0$$

$$\beta_i = b_i - \sum_{l=1}^i a_l\beta_{n-1-l} \quad 1 \leq i \leq n-1$$

then

$$\beta_0 = b_0 = 0 \quad \beta_1 = b_1 - a_1\beta_0 = \frac{b}{m} \quad \beta_2 = \frac{k}{m} - \left(\frac{b}{m}\right)^2$$

From Eq. 255, recall

$$x_1 = y - \beta_0u$$

$$x_i = \frac{d^{i-1}y}{dt^{i-1}} - \sum_{j=0}^{i-2} \beta_j \frac{d^{i-1-j}u}{dt^{i-1-j}} - \beta_{i-1}u = \dot{x}_{i-1} - \beta_{i-1}u \quad 2 \leq i \leq n$$

$$\begin{aligned} x_1 &:= y - b_0u = y \\ \implies x_2 &:= \dot{x}_1 - \frac{b}{m}u \end{aligned}$$



From Eq. 257

$$\begin{aligned}\dot{x}_i &= x_{i+1} + \beta_i u \quad 1 \leq i \leq n-1 \\ \dot{x}_n &= -\sum_{l=1}^n a_{n+1-l} x_l + \beta_n u\end{aligned}$$

and so

$$\begin{aligned}\Rightarrow \quad \dot{x}_1 &= x_2 + \frac{b}{m} u \\ \dot{x}_2 &= -\frac{k}{m} x_1 - \frac{b}{m} x_2 + \left( \frac{k}{m} - \left( \frac{b}{m} \right)^2 \right) u\end{aligned}$$

So output equation  $y = x_1$ ,

(261)

$$\dot{\mathbf{x}} = \begin{bmatrix} 0 & 1 \\ \frac{-k}{m} & \frac{-b}{m} \end{bmatrix} \mathbf{x} + \begin{bmatrix} \frac{b}{m} \\ \frac{k}{m} - \left( \frac{b}{m} \right)^2 \end{bmatrix} u$$

**Example 3-5.** An inverted pendulum mounted on motor-driven cart. This is a model of the attitude control of a space booster on takeoff.

44. DISCRETE CONTROL, S-DOMAIN TO Z-DOMAIN

cf. [Discrete control 1: Introduction and overview](#) Douglas (2017) [46]  
Recall Laplace transform of unit step function is  $\frac{1}{s}$  (cf. pp. 161, Ogata (2009) [47]).  
Impulse trains aren’t usually the discrete input when you have a discrete system that interacts with a continuous system.

**44.1. Zero-order hold (ZOH) or step-invariant method.** Discrete command or measurement  $\rightarrow$  hold constant  $\rightarrow$  next sample time  $\rightarrow$  hold constant

We want our discretized system to incorporate the effects of the zero order hold. Let’s walk through the math.  
We need to start with some definitions.

(262)

$$\begin{aligned}z &= e^{Ts} \\ T &= \text{sample period} \\ s &\in \mathbb{C}\end{aligned}$$

(263)

$$z - \text{transform} \rightarrow \mathcal{Z}(f[n]) = \sum_{k=0}^{\infty} f[k] z^{-k}$$

Now we can look at the zero order hold.  
 $v_k := v[k] \xrightarrow{\text{ZOH}}$  step functions  $v(t)$ . We can look at the  $k$ th pulse, between  $kT$  and  $(k+1)T$ , height is the value of  $v_k[k]$ .  
 $k$ th pulse can be created with 2 step functions

$$\begin{aligned}&v_k[k] \cdot U(t - kT) \\ &+ -v_k[k] \cdot U(t - (k+1)T)\end{aligned}$$

where  $U(t - kT) = \begin{cases} 1 & \text{if } t \geq kT \\ 0 & \text{if } t < kT \end{cases}$

So that the single pulse at index  $k \implies$

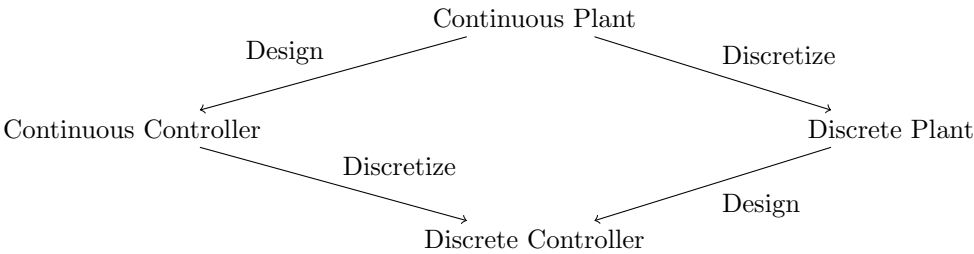
(264)

$$v_k[k][(U(t - kT) - U(t - (k+1)T))]$$

cf. [Discrete control 3: Designing for the zero-order hold](#), Douglas (2017) [46].  
Digital computer supplies sequence of discrete voltages.

Then DAC, which is a physical piece of electrical hardware, holds the voltage constant between commands and then steps up or down to new voltage of the next sample time.

Recall



You can use any discretization method to get **a** discrete transfer function.  
But it’s misleading to imply we’ll use the ZOH (Zero-Order Hold) method to get the discrete controller we **want!**  
cf. [Discrete control 6: z-plane warping and the bilinear transform](#), Douglas (2017) [46].  
We have a large amount of theory available to us (for continuous systems) that’ll help us craft the perfect analog filter.  
Let’s try our friend’s suggestion on simple notch filter.  
Design requirements:  $\omega_0 = 0.1 \frac{\text{rad}}{\text{sec}} \rightarrow$  critical frequency.  
 $Q = 1$

Given

$$G(s) = \frac{s^2 + 0.01}{s^2 + 0.1s + 0.01}$$

Let’s try the bilinear transform on this in Matlab.

(265)

$$\omega_d = \frac{2}{T} \arctan \left( \omega_a \frac{T}{2} \right)$$

Linear at low frequencies compared to Nyquist.  
So what can we do?  
Prewarp  $\omega_a$ !  
Apply tan equation to  $\omega_a$  before converting with the bilinear transform, so that

$$s \mapsto \frac{\omega_a}{\tan \left( \frac{\omega_a T}{2} \right)} \frac{z - 1}{z + 1}$$

45. SPACECRAFT ATTITUDE CONTROL

**45.1. Reaction Wheels.** cf. Slide 21/45, M. Peet, Lecture 15.

A momentum exchange device applies torque to wheel to spin up wheel.  
- An equal and opposite amount of torque is imparted to the spacecraft.  
The resulting angular momentum of wheel and craft each are equal in magnitude and opposite in direction by angular momentum conservation.

Consider rotation about x-axis.  
Let  $J_x$  be moment of inertia of spacecraft about x-axis.  
Let  $I_x$  be moment of inertia of flywheel.  
By anglar umomentum conservation,

$$I_x(\omega_f + \omega_s) + J_x \omega_s = 0$$

where  $\omega_s \equiv$  angular velocity of craft in inertial space.  
 $\omega_f \equiv$  angular velocity of flywheel with respect to craft.

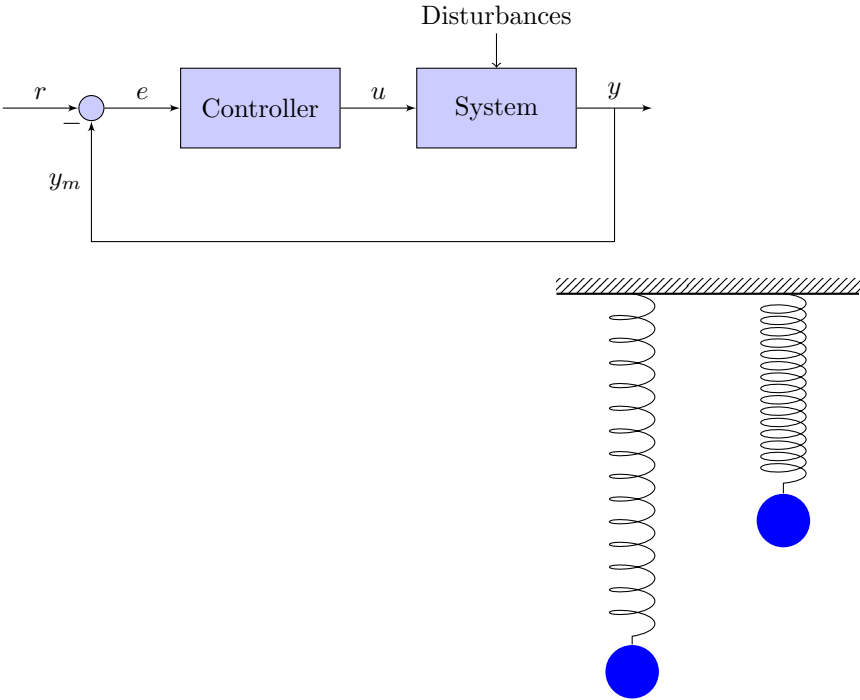
So if craft has some angular velocity  $\omega_s$  in  $\hat{b}_1$ -direction (body-frame axis 1?) and reaction wheel is aligned with this axis, null out velocity by spinning up to

46. MORE RESOURCES

Daniel Block, Karl Åström, Mark W. Spong. **The Reaction Wheel Pendulum**, i.e. Block, Åström, and Spong (2007) [48]  
Kuo. **Digital Control Systems**.  
Sidi. **Spacecraft Dynamics and Control**.

46.1. **Miscellaneous.** barfoot wedge  $SO(3)$ . barfoot lie algebra. [http://asrl.utias.utoronto.ca/~tdb/bib/barfoot\\_ser17.pdf](http://asrl.utias.utoronto.ca/~tdb/bib/barfoot_ser17.pdf)

Simple example Latex diagram of control system:



Part 15. Orbital Mechanics

Curtis (2020) [41]  
See pp. 54, Ch. 1 "Dynamics of Point Masses" of Curtis (2020) [41].

**Problem 1.23.** Use an RK solver to solve the nonlinear Lorenz equations, due to the American meteorologist and mathematician E.N. Lorenz (1917-2008):

$$\begin{aligned}\dot{x} &= \sigma(y - x) \\ \dot{y} &= x(\rho - z) - y \\ \dot{z} &= xy - \beta z\end{aligned}$$

Start off by using the values Lorenz (1963) used in his paper (namely,  $\sigma = 10, \beta = 8/3$ , and  $\rho = 28$ ). For initial conditions use  $x = 0, y = 1$ , and  $z = 0$  at  $t = 0$ .

So from "**Deterministic Nonperiodic Flow**"

cf. Wie (2008) [42]

47. N-BODY PROBLEM

cf. *N*-body problem, Sec. 1.2, Bate, Mueller, and White (1971) [35]  
Force  $\mathbf{F}_{gn}$  exerted on  $m_i$  by  $m_n$  is  $\mathbf{F}_{gn} = \frac{-Gm_i m_n}{(r_{ni})^3} \mathbf{r}_{ni}$  (cf. Eq. (1.2-1), Bate, Mueller, and White (1971) [35]) where  $\mathbf{r}_{ni} = \mathbf{r}_i - \mathbf{r}_n$   
Force  $\mathbf{F}_g$  of all gravitational forces acting on  $i$ th body.

(266)

$$\mathbf{F}_g = -Gm_i \sum_{\substack{j=1 \\ j \neq i}}^n \frac{m_j}{(r_{ji})^3} \mathbf{r}_{ji} \text{ where } \mathbf{r}_{ji} \equiv \mathbf{r}_i - \mathbf{r}_j$$

cf. Eq. (1.2-5) Bate, Mueller, and White (1971) [35]  
Consider  $\mathbf{F}_{\text{tot}} = \frac{d\mathbf{p}}{dt} = \frac{d}{dt}(m\mathbf{v}) = \dot{m}\mathbf{v} + m\frac{d\mathbf{v}}{dt}$  or

(267)

$$\frac{d\mathbf{v}}{dt} = \frac{\mathbf{F}_{\text{tot}}}{m} - \frac{\dot{m}}{m} \mathbf{v}$$

cf. Eq. (1.2-9), Bate, Mueller, and White (1971) [35].  
Assume  $m_2$  is an earth satellite,  $m_1$  is the earth; remaining masses  $m_3, m_4, \dots m_n$  may be moon, sun, planets. Assume drag and other external forces aren't present. The only remaining forces are gravitational.

$$\begin{aligned}\frac{d\mathbf{v}_1}{dt} &= -G \sum_{j=2}^n \frac{m_j}{(r_{j1})^3} \mathbf{r}_{j1} \\ \frac{d\mathbf{v}_2}{dt} &= -G \sum_{\substack{j=1 \\ j \neq 2}}^n \frac{m_j}{(r_{j2})^3} \mathbf{r}_{j2}\end{aligned}$$

Consider  $\mathbf{r}_{12} = \mathbf{r}_2 - \mathbf{r}_1$ , so  $\ddot{\hat{r}}_{12} = \ddot{\mathbf{r}}_2 - \ddot{\mathbf{r}}_1$ , so

(268)

$$\begin{aligned}\frac{d\mathbf{v}_{12}}{dt} &= -G \left[ \sum_{j=3}^n \left( \frac{m_j}{(r_{j2})^3} \mathbf{r}_{j2} - \frac{m_j}{(r_{j1})^3} \mathbf{r}_{j1} \right) + \frac{m_1}{(r_{12})^3} \mathbf{r}_{12} - \frac{m_2}{(r_{21})^3} \mathbf{r}_{21} \right] = -G \left[ \sum_{j=3}^n (m_j) \left( \frac{\mathbf{r}_{j2}}{(r_{j2})^3} - \frac{\mathbf{r}_{j1}}{(r_{j1})^3} \right) + (m_1 + m_2) \frac{\mathbf{r}_{12}}{(r_{12})^3} \right] = \\ &= -G \frac{(m_1 + m_2) \mathbf{r}_{12}}{(r_{12})^3} - \sum_{j=3}^n G m_j \left( \frac{\mathbf{r}_{j2}}{(r_{j2})^3} - \frac{\mathbf{r}_{j1}}{(r_{j1})^3} \right)\end{aligned}$$

since  $\mathbf{r}_{12} = -\mathbf{r}_{21}$ . cf. Eq. (1.2-17) Bate, Mueller, and White (1971) [35]  
 $m_2$  is the mass of a satellite,  $m_1$  is the mass of the earth.  $\ddot{\mathbf{r}}_{12} = \dot{\mathbf{v}}_{12}$  is the acceleration of the satellite relative to Earth.  
Effect of the last term is to account for the perturbing effects of the moon, sun, planets on near earth satellite.

47.1. **2-body problem, equation of relative motion.** cf. 2 body Problem. 1.3.2 The Equation of Relative Motion. Bate, Mueller, and White (1971) [35].

Consider 2 assumptions:

- (1) The bodies are spherically symmetric. This enables us to treat bodies as though their masses were concentrated at their centers.
- (2) There are no external nor internal forces acting on system other than gravitational forces acting along line joining centers of the 2 bodies.

For masses  $m, M$ , define

(269)

$$\mathbf{r} := \mathbf{r}_m - \mathbf{r}_M$$

Apply Newton's laws in inertial frame:

$$(270) \quad \begin{aligned} m\ddot{\mathbf{r}}_m &= -\frac{GMm}{r^2} \frac{\mathbf{r}}{r} & \ddot{\mathbf{r}}_m &= -\frac{GM}{r^2} \frac{\mathbf{r}}{r} \\ M\ddot{\mathbf{r}}_M &= \frac{GMm}{r^2} \frac{\mathbf{r}}{r} & \ddot{\mathbf{r}}_M &= \frac{Gm}{r^2} \frac{\mathbf{r}}{r} \\ && \implies \ddot{\mathbf{r}} &= \frac{-G(M+m)}{r^3} \mathbf{r} \end{aligned}$$

cf. (1.3-1), (1.3-2) and (1.3-3) Bate, Mueller, and White (1971) [35].

If  $M \gg m$ , then  $G(M+m) \approx GM$ . Define  $\mu \equiv GM$ , so that

$$(271) \quad \ddot{\mathbf{r}} = \frac{-\mu}{r^3} \mathbf{r}$$

cf. Eq. (1.3-4) Bate, Mueller, and White (1971) [35]

47.1.1. *Mechanical energy conservation.* cf. 1.4.1 "Conservation of Mechanical Energy", pp. 15, Bate, Mueller, and White (1971) [35]

Dot multiply  $\ddot{\mathbf{r}} + \frac{\mu}{r^3} \mathbf{r} = 0$  Eq. 271, i.e. Eq. (1.3-4) of Bate, Mueller, and White (1971) [35] by  $\dot{\mathbf{r}}$ :

$$\dot{\mathbf{r}} \cdot \ddot{\mathbf{r}} + \frac{\mu}{r^3} \dot{\mathbf{r}} \cdot \mathbf{r} = \mathbf{v} \cdot \dot{\mathbf{v}} + \frac{\mu}{r^3} \mathbf{r} \cdot \dot{\mathbf{r}} = 0$$

Now note that

$$\frac{d}{dt} \left( \frac{1}{\sqrt{x_1^2 + \cdots + x_N^2}} \right) = \frac{-1}{2} (x_1^2 + \cdots + x_N^2)^{-3/2} (2x_1 \dot{x}_1 + \cdots + 2x_N \dot{x}_N)$$

so then

$$\begin{aligned} \frac{1}{2} \frac{d}{dt} (v^2) &= \frac{1}{2} \frac{d}{dt} (v_i v_i) = v_i \dot{v}_i = \mathbf{v} \cdot \dot{\mathbf{v}} & \frac{d}{dt} \left( \frac{-\mu}{r} \right) &= \frac{\mu}{r^3} (r_i \dot{r}_i) \\ &\implies \frac{d}{dt} \left( \frac{v^2}{2} + \frac{-\mu}{r} + C \right) = 0 \end{aligned}$$

So we obtain mechanical energy conservation. Let  $c = 0$ , setting the zero of the potential energy at infinity, effectively. Then

$$(272) \quad E = \frac{v^2}{2} - \frac{\mu}{r}$$

cf. Eqn. (1.4-2) Bate, Mueller, and White (1971) [35].

47.1.2. *Angular Momentum conservation.* Cross multiply Eq. 271, i.e. Eqn. (1.3-4) by  $\mathbf{r}$

$$\mathbf{r} \times \ddot{\mathbf{r}} + \mathbf{r} \times \frac{\mu}{r^3} \mathbf{r} = 0 = \mathbf{r} \times \ddot{\mathbf{r}} = \frac{d}{dt} (\mathbf{r} \times \dot{\mathbf{r}}) = \frac{d}{dt} (\mathbf{r} \times \mathbf{v})$$

Bate, Mueller, and White (1971) [35] uses this notation:

$$(273) \quad \mathbf{h} = \mathbf{r} \times \mathbf{v}$$

in Eq. (1.4-3) of Bate, Mueller, and White (1971) [35], where  $\mathbf{h}$  is called the specific angular momentum. Instead, let's use my notation:

$$(274) \quad \mathbf{l} = \mathbf{r} \times \mathbf{v}$$

So angular momentum conservation says

$$(275) \quad \frac{d\mathbf{l}}{dt} = 0$$

Bate, Mueller, and White (1971) [35] defines

"up"  $\equiv$  away from the center of the earth,

"down"  $\equiv$  toward the center of the earth.

So local vertical at location of satellite coincides with direction of vector  $\mathbf{r}$ .

Define direction of  $\mathbf{v}$  by specifying angle  $\gamma$  that  $\mathbf{v}$  makes with local vertical  $\equiv$  zenith angle.

Define  $\phi \equiv$  flight-path elevation angle  $\equiv$  "flight-path angle"  $\equiv$  angle between  $\mathbf{v}$  and local horizontal plane.

$$\implies |\mathbf{l}| = rv \sin \gamma$$

Since  $\gamma, \phi$  complementary angles (i.e.  $\frac{\pi}{2} - \phi = \gamma$ ) then

$$(276) \quad \boxed{|\mathbf{l}| = rv \cos \phi}$$

cf. (1.4-4) of Bate, Mueller, and White (1971) [35].

47.1.3. *Trajectory Equation for 2-body problem.* cf. Sec. 1.5. "The Trajectory Equation" of Bate, Mueller, and White (1971) [35].

Recall again Eq. 271

$$\ddot{\mathbf{r}} = -\frac{\mu}{r^3} \mathbf{r}$$

Apply the cross product by  $\mathbf{l}$ :

$$\ddot{\mathbf{r}} \times \mathbf{l} = \frac{-\mu}{r^3} \mathbf{r} \times \mathbf{l} = \frac{-\mu}{r^3} [(\mathbf{r} \cdot \mathbf{v}) \mathbf{r} - r^2 \mathbf{v}] = \mu \frac{d}{dt} \left( \frac{\mathbf{r}}{r} \right)$$

since

$$\frac{d}{dt} \left( \frac{\mathbf{r}}{r} \right) = \frac{-1}{r^3} (r_i \dot{r}_i) \mathbf{r} + \frac{1}{r} \mathbf{v} = \frac{-\mathbf{r} \cdot \mathbf{v} \mathbf{v}}{r^3} + \frac{\mathbf{v}}{r}$$

Consider that  $\dot{\mathbf{r}} \times \mathbf{l} = \mathbf{v} \times \mathbf{l} = \epsilon_{ijk} v_j l_k \mathbf{e}_i$ . Then

$$\frac{d}{dt} (\mathbf{v} \times \mathbf{l}) = \epsilon_{ijk} \mathbf{e}_i \dot{v}_j l_k + \epsilon_{ijk} \mathbf{e}_i v_j \dot{l}_k = \ddot{\mathbf{r}} \times \mathbf{l}$$

where we had used angular momentum conservation from Eq. 275:  $\frac{d\mathbf{l}}{dt} = 0$ .

So then

$$\frac{d}{dt} (\mathbf{v} \times \mathbf{l}) = \mu \frac{d}{dt} \left( \frac{\mathbf{r}}{r} \right) \xrightarrow{f} \mathbf{v} \times \mathbf{l} = \mu \frac{\mathbf{r}}{r} + \mathbf{B} \xrightarrow{r} \mathbf{r} \cdot (\mathbf{v} \times \mathbf{l}) = \mu r + \mathbf{B} \cdot \mathbf{r}$$

where  $\mathbf{B}$  is an integration constant.

Now  $\mathbf{a} \cdot (\mathbf{b} \times \mathbf{c}) = (\mathbf{a} \times \mathbf{b}) \cdot \mathbf{c}$  since  $a_i \epsilon_{ijk} b_j c_k = (\epsilon_{ijk} a_i b_j) c_k$ , so then

$$(277) \quad l^2 = \mu r + \mathbf{B} \cdot \mathbf{r} = \mu r + Br \cos \nu \text{ or } r = \frac{l^2/\mu}{1 + \frac{B}{\mu} \cos \nu} \equiv \frac{l^2/\mu}{1 + \frac{B}{\mu} \cos \varphi}$$

cf. pp. 24 "Two-Body Orbital Mechanics" of Bate, Mueller, and White (1971) [35]

Bate, Mueller, and White (1971) [35] describes the width of each (conic section) curve at the focus as a positive dimension called the *latus rectum* and labels it  $2p$ . It is equated as Eqn. (1.5-6) in Bate, Mueller, and White (1971) [35]:

$$(278) \quad p = a(1 - e^2)$$

Furthermore,

periapsis - pt. nearest prime focus

apoapsis - pt. farthest from the prime focus

The distance from the prime focus to either the periapsis or apoapsis (where it exists; for open curves, parabolas, and hyperbolas, an apoapsis has no physical meaning) are given from  $r = \frac{a(1-e^2)}{1+e \cos \theta}$  (polar equation of a conic section with focus at origin):

$$(279) \quad r_{\min} = r_{\text{periapsis}} = \frac{p}{1 + e \cos 0} = \frac{p}{1 + e} = r_p$$

cf. Eq. (1.5-7) of Bate, Mueller, and White (1971) [35]

$$(280) \qquad r_{\max} = r_{\text{apoapsis}} = \frac{p}{1 + e \cos \pi} = \frac{p}{1 - e} = r_a$$

cf. Eq. (1.5-8) of Bate, Mueller, and White (1971) [35]

For an ellipse, with  $e < 1$ ,  $r_p = \frac{a(1-e^2)}{1+e}$ ;  $r_a = \frac{a(1-e^2)}{1-e} = a(1 + e)$

For a hyperbola,  $e > 1$ ,  $r_p = -a(e - 1)(e + 1)/(e + 1) = -a(e - 1)$ .

47.1.4. *Bate, Mueller, and White's so-called "Eccentricity Vector"*. Recall the derivation of Eqn. (1.5-3) of Bate, Mueller, and White (1971) [35],  $r = \frac{h^2/\mu}{1+(B/\mu)\cos \nu}$ . Now recall from Eq. ??:

$$\mathbf{v} \times \mathbf{l} = \frac{\mu \mathbf{r}}{r} + \mathbf{B} \implies \mathbf{B} = \mathbf{v} \times \mathbf{l} - \frac{\mu \mathbf{r}}{r}$$

Bate, Mueller, and White (1971) [35] defines the eccentricity vector as

$$\mathbf{e} = \frac{\mathbf{v} \times \mathbf{l}}{\mu} - \frac{\mathbf{r}}{r}$$

Compare this against the Laplace-Runge-Lenz vector of Sec. 3.9 "The Laplace-Runge-Lenz Vector" in Goldstein, et. al. [11]. cf. 1.6 "Relating  $\mathbf{E}$  and  $\mathbf{l} \equiv \mathbf{h}$  to the Geometry of an Orbit in Bate, Mueller, and White (1971) [35].

Bate, Mueller, and White (1971) [35] compares  $r = \frac{l^2/\mu}{1+B/\mu\cos \nu}$  to  $r = \frac{p}{1+e\cos \nu}$ , Eqns. (1.5.3), (1.5.4), respectively, from Bate, Mueller, and White (1971) [35]. To obtain

$$(281) \qquad p = l^2/mu$$

cf. Eqn. (1.6-1) of Bate, Mueller, and White (1971) [35].

Since  $\mathbf{l}$  is conserved, then

$$(282) \qquad l = r_p v_p = r_a v_a$$

cf. (1.6-2) of Bate, Mueller, and White (1971) [35], i.e. at the periapsis or apoapsis of any conic orbit the velocity vector (which is always tangent to the orbit) is directed horizontally and flight-path angle  $\nu$  is 0.

From Eqns. (1.4-2), (1.5-7), (1.5-6), (1.6.1) of Bate, Mueller, and White (1971) [35],

$$\mathcal{E} = \frac{v^2}{2} - \frac{\mu}{r} = \frac{l^2}{2r_p^2} - \frac{\mu}{r_p} \xrightarrow{r_p=a(1-e)} \frac{l^2 - 2\mu a(1 - e)}{2r^2p} \xrightarrow{l^2=\mu a(1-e^2)} \frac{-\mu}{2a}$$

Since  $l$  determines  $p = a(1 - e^2)$  so  $e = \sqrt{1 - \frac{p}{a}}$  and so  $p = l^2/\mu$  and  $a = -\frac{\mu}{2\mathcal{E}}$ . Then

$$(283) \qquad e = \sqrt{1 + \frac{2\mathcal{E}l^2}{\mu^2}}$$

cf. Eq (1.6-4) of Bate, Mueller, and White (1971) [35].

#### 48. ORBIT DETERMINATION FROM OBSERVATIONS; COORDINATE SYSTEMS

cf. Ch. 2 "Orbit Determination From Observations", Bate, Mueller, and White (1971) [35].

48.1. **Coordinate Systems.** cf. Sec. 2.2 "Coordinate Systems", Bate, Mueller, and White (1971) [35].

To describe coordinate systems, we'll give position of origin, orientation of the fundamental plane ( $X - Y$  plane), principal direction (i.e. direction of  $X$ -axis), and direction of  $Z$ -axis.

$Z$ -axis must be perpendicular to fundamental plane.

48.1.1. *Heliocentric-Ecliptic Coordinate System.* cf. Sec. 2.2.1 "The Heliocentric-Ecliptic Coordinate System", Bate, Mueller, and White (1971) [35]

heliocentric-ecliptic system has origin at sun center. So, origin at sun center.

$X_\epsilon - Y_\epsilon$  or fundamental plane coincides with "ecliptic" which is plane of earth's revolution around sun.

line-of-intersection of ecliptic plane and earth's equatorial plane defines direction of  $X_\epsilon$ -axis.

On first day of spring, line joining earth's center and sun's center points in direction of positive  $X_\epsilon$ -axis; called *vernal equinox direction*.

Earth axis of rotation shifts slowly over centuries, i.e. precession, so line-of-intersection of earth's equator and ecliptic shifts slowly.

As a result, heliocentric-ecliptic system isn't really an inertial reference frame and where precession is required, it'd be necessary to specify  $XYZ_\epsilon$  coordinates of an object were based on vernal equinox direction of a particular year or "epoch."

48.1.2. *Geocentric-Equatorial Coordinate System.* geocentric-equatorial system has origin at earth's center.

#### Part 16. GPS, Geodesy

D. K. Olson, "Converting Earth-Centered, Earth-Fixed Coordinates to Geodetic Coordinates" IEEE Transactions on Aerospace and Electronic Systems, 32 (1996) 473-476 Olson algorithm

#### 49. GEODESY

Mike Craymer (2023)'s Geodetic Toolbox in Matlab.

#### Part 17. Real-time systems

#### 50. PARTIAL ORDERING OF EVENTS IN A DISTRIBUTED SYSTEM

cf. pp. 558, Communications of the ACM, July 1978, Vol. 21, Number 7, Lamport (1978) [49]

Assume system composed of a collection of processes. Each process consists of a sequence of events.

A single process is defined to be a set of events with an a priori total ordering.

Recall for set  $P$ , ordering  $\leq$ , then  $(P, \leq)$  **partially ordered** if

- (1) reflexivity:  $a \leq a$
- (2) anti-symmetry:  $(a \leq b \text{ and } b \leq a) \rightarrow a = b$
- (3) transitive:  $a \leq b \text{ and } b \leq c \implies a \leq c$

We can extend our definition to allow a process to split into distinct subprocesses (EY: if needed).

Assume sending or receiving a message is an event in a process.

define "happened before" relation:

**Definition 10** ("happened before" relation). *Given set of events  $E$  of a system, with set of event  $P \subseteq E$  s.t.  $\forall$  event  $\in P$  is on the same process  $P$ , relation  $<_p$  (Lamport's notation is " $\rightarrow$ ") is the "smallest" relation s.t.*

- (1) *irreflexivity:  $a \not\prec_p a \quad \forall a \in E$  (interpret: event cannot happen before itself)*
- (2) *concurrency: if  $a \not\prec_p b$  and  $b \not\prec_p a$ , then  $a, b$  **concurrent**. (contrapositive to concurrency)  $a, b$  not concurrent if  $a <_p b$  or  $b <_p a$*

*EY: compare this with anti-symmetry, if  $a \leq_p b$  and  $b \leq_p a$ , then  $a =_p b$ ; this conclusion requires from concurrency that if  $a \not\prec_p b$ , then  $b \leq_p a$ . Then consider the extra claim that  $a \leq_p b$  is equivalent to  $a <_p b$  or  $a =_p b$  (concurrent).*

(3) *transitivity: if  $a <_p b$  and  $b <_p c$ , then  $a <_p c$*

*Can be shown that if  $a \leq_p b$  and  $b \leq_p c$ , then  $a \leq_p c$  (consider all permutations/cases: if  $a =_p b$ ,  $b =_p c$ , then  $a =_p c$ , if  $a =_p b$ ,  $b <_p c$ , then  $a <_p c$ , ...)*

*Intepret:  $\forall a, b \in P$  (interpret  $a, b$  are events in the same process),  $a <_p b$  means  $a$  comes before  $b$ .*

*Interpret, if processes  $P, Q \subset E$ , and event  $a$  is sending of a message by process  $P$ , and event  $b$  is receipt of same message by another process  $Q$ , then  $a <_p b$ .*

EY: if  $a \not\leq_p b$ , equivalent to  $b \leq_p a$ ,  $\forall a, b \in E$ , then  $a \leq_p a$ ,  $\forall a \in E$ , if  $a \leq_p b$  and  $b \leq_p a$ , then  $a =_p b$ , and  $a \leq_p b$  and  $b \leq_p c$ , then  $a \leq_p c$ .

So if  $a \not\leq_p b$  equivalent to  $b \leq_p a$ , and  $a =_p b$  is equivalent to concurrency,  $<_p$  **implies partial ordering**  $\leq_p$ .

Also consider for  $<_p$  in Def. 10 to be interpreted in terms of causation, and two events are concurrent if neither can causally affect the other.

## 51. LOGICAL CLOCKS

cf. pp. 559, Communications of the ACM, July 1978, Vol. 21, Number 7, Lamport (1978) [49]

Introduce clocks in the system.

Begin with abstract point of view such that clock is just a way of assigning number to an event, where number is thought of as time at which event occurred.

Define  $\forall$  process  $P_i$ , a clock  $C_i$ ,

$$(284) \quad \begin{aligned} C_i : P_i &\rightarrow \mathbb{R} \text{ or } \mathbb{Z} \\ C_i : P_i &\rightarrow \mathbb{Z} \text{ for system clock or digital computer clock} \end{aligned}$$

For clock  $i$  on process  $i$ , i.e. clock  $C_i$  for each process  $P_i$ , then  $\forall$  event  $a \in P_i$ ,

$$(285) \quad C_i \langle a \rangle \in \mathbb{Z}$$

Entire system of clocks represented by mapping  $C$  s.t.

$$(286) \quad \begin{aligned} C : b &\mapsto C \langle b \rangle \\ C \langle b \rangle &= C_j \langle b \rangle \text{ if } b \text{ is an event in process } P_j \end{aligned}$$

(”For now”) we make no assumption about the relation of the numbers  $C_i \langle a \rangle$  to physical time, so we can think of clocks  $C_i$  as logical rather than physical clocks. They maybe implemented by counters with no actual timing mechanism.

Base definition on **order in which events occur**; we cannot base our definition of correctness on physical time, since that would require introducing clocks which keep physical time.

Lamport makes this strongest reasonable condition: if an event  $a$  occurs before another event  $b$ , then  $a$  should happen at an earlier time than  $b$ :

**Proposition 1** (Clock Condition).  $\forall$  events  $a, b \in E$ ,  
if  $a <_p b$ , then  $C \langle a \rangle < C \langle b \rangle$

Lamport says that we can’t expect converse condition, since that’d imply any 2 concurrent events must occur at the same time:

(converse to Prop. 1) if  $a \not\leq_p b$ , then  $C \langle a \rangle \geq C \langle b \rangle$  or if  $a =_p b$  or  $a >_p b$ , then  $C \langle a \rangle \geq C \langle b \rangle$ . Further, if  $a \not\leq_p b$ , then this implies  $C \langle a \rangle = C \langle b \rangle$  (since it cannot be that  $C \langle a \rangle > C \langle b \rangle$ , we eliminated this possibility out of the two implied by  $a \geq_p b$ ). But we *don’t want* for 2 concurrent events to necessarily have equal times (because who’s to say that they must occur at the same time).

e.g. On process  $P$ ,  $p_2 <_p p_3$ , and  $C_p(p_2) < C_p(p_3)$ , because events on the same process are totally ordered.

Given,

$$p_1 <_p q_2$$

$$q_1 <_p p_2$$

$$q_2 <_p q_3$$

$p_2 \not\leq_p q_2$  and  $q_2 \not\leq_p p_2$ , and since  $q_2 < q_3$ , then

since  $p_2 \not\rightarrow q_3$  and  $q_2 \not\rightarrow p_3$ ,  $p_2, q_3$  concurrent.

since  $p_3 \not\rightarrow q_3$  and  $q_3 \not\rightarrow p_3$ ,  $p_3, q_3$  concurrent.

Clock condition is satisfied if

**Proposition 2** (Clock conditions 1, 2). C1 If  $a, b \in$  process  $P_i$ , and  $a < b$ ;  $C_i \langle a \rangle < C_i \langle b \rangle$

C2 If  $a$  is sending a message by process  $P_i$  and  $b$  is receipt of that message by process  $P_j$ , then  $C_i \langle a \rangle < C_j \langle b \rangle$

cf. pp, 560, Lamport (1978) [49]

Let’s assume that processes are algorithms, and events represent certain actions during their execution. Process  $P_i$ ’s clock is represented by  $C_i$  during event  $a$ .

For condition C1, if  $a, b \in P_i$ ,  $a <_p b$ , then  $C_i(a) < C_i(b)$ , so then processes need to obey the following

**Proposition 3** (Implementation Rule 1 (IR1)). Each process  $P_i$  increments  $C_i$  between any 2 successive events.

EY (20210408) i.e.  $\forall$  process  $P_i$ , let sequence of events  $a_k \in P_i$  s.t.  $\forall a_k$  totally ordered on  $P_i$ , i.e.  $a_k <_p a_{k+1}$ ,  $\forall k$ , then

$$(287) \quad \boxed{C_i(a_k) < C_i(a_{k+1}) = C_i(a_k) + 1}$$

To meet condition C2, if  $a$  sends message by process  $P_i$ ,  $b$  receives same messsage by process  $P_j$ , then

$$C_i(a) < C_j(b)$$

We require that each message  $m$  contain a timestamp  $T_m$  which equals the time at which message was sent.

Upon receiving a message timestamped  $T_m$ , a process must advance its clock to be later than  $T_m$ , More precisely,

**Proposition 4** (Implementation Rule 2 (IR 2)). (1) If event  $a$  sends message  $m$  by process  $P_i$ , then message  $m$  contains timestamp  $T_m = C_i \langle a \rangle$

(2) Upon receiving a message  $m$ , process  $P_j$  sets  $C_j$  greater than or equal to its present value and greater than  $T_m$

EY (20210408), i.e.

If event  $a \in P_i$  sends message  $m_a$  contains timestamp  $T_a = C_i(a)$

Upon receiving message  $m_a$ , event  $r_a$ , process  $P_j$  sets

$$(288) \quad \boxed{C_j(r_a) = \max(C_j(r_a), T_a) + 1}$$

The addition, increment, of 1 makes sense since the message can never be received before or at the same time of sending it (cf. *Steven Van Dorpe*)

(remember that  $C_j(r_a) = C_j(q) + 1$  where  $q$  is event immediately before  $r_a$ ,  $q <_p r_a$  (and no other event))

IR2 insures C2 satisfied; hence IR1, IR2 imply Clock Condition satisfied.

51.1. **Ordering Events Totally.** Define relation  $\implies$ .

Use any arbitrary total ordering  $<_{(P)}$  of the processes to break ties.

If event  $a$  in process  $P_i$ , event  $b$  in process  $P_j$ , then  $a \implies b$  iff either

- (1)  $C_i \langle a \rangle < C_j \langle b \rangle$  or
- (2)  $C_i \langle a \rangle < C_j \langle b \rangle$  and  $P_i < P_j$

ordering  $\implies$  depends upon system of clocks  $C_i$  and is not unique.

Example: use total ordering of events solving the following version of the *mutual exclusion problem* (EY: mutex?).

Consider system composed of fixed collection of processes  $P_i$ , (e.g.  $N = 1$ ,  $P$ , e.g.  $N = 2$ ,  $P_i$ ,  $P_j$ ). Processes share a *single* resource (resource is external, not part of any of the processes, to any of the processes, I believe).

We wish to find an algorithm satisfying following 3 conditions:

- (I) A process  $P_i$  which has been granted *the* resource must release it before it can be granted to another process  $P_j$
- (II) Different requests for *the* resource must be granted in order in which they’re made
- (III) If every process which is granted the resource eventually releases it, then every request is eventually granted

Implement:

Each process maintains own request queue which is never seen by any other process.



- (1) To request resource, process  $P_i$ , sends  $T_m : P_i$  "requests resource" to every other process (e.g. if  $N = 2$ , to  $P_{i+1}$ ; if  $N = 1$ , no send) *and* puts *that message*  $T_m : P_i$  "requests resource" on its request queue
- (2) When process  $P_j$  receives message  $T_m : P_i$  "requests resource", it places it on its request queue *and* sends (timestamped) acknowledgement message to  $P_i$  ( $N = 1$  none,  $N = 2$ ,  $P_j$  receives and sends ack to  $P_i$ )
- (3) To release resource  $P_i$  removes *any*  $T_m : P_i$  "requests resource" message from its request queue, *and* sends (timestamped)  $P_i$  "releases resource" message to every other process ( $N = 1$ , removes all  $T_{m'} : P_i$  "requests resource" messages;  $N = 2$ ,  $P_i$  removes any  $T_m : P_i$  "requests resources" and  $P_j$  gets sent  $T_{m'} : P_i$  "releases resource")
- (4) When  $P_j$  receives  $P_i$  "releases resource", removes any  $T_m : P_i$  "requests resource" message from its request queue.
- (5)  $P_i$  granted resource if
  - (i)  $\exists T_m : P_i$  "requests resource" message in its request queue ordered before any *other request*
  - (ii)  $P_i$  received message from  $\forall P_{i'}$  timestamped later than  $T_m$  ( $N = 1$ , none;  $N = 2$ ,  $P_i$  received ack (acknowledgement) from  $P_j$  at  $T_{m'} > T_m$ )

5(ii)  $\rightarrow$  guarantees  $P_i$  has learned about all requests which preceded its current request

rules 3,4 are only ones which delete messages from request queue  $\rightarrow$  Condition I

Condition II from total ordering  $\implies$  extends partial ordering

Rule 2  $\implies$  guarantees after  $P_i$  requests resource, rule 5(ii) holds

Rules 3,4  $\implies$  if  $\forall P_i$  granted resource eventually releases it, then rule 5(i),  $\implies$  Condition III

$\implies$

$\mathbf{C}$  := set of possible commands, e.g.  $\mathbf{C} = \{ \text{rule 1, rule 3} \} =$  all commands  $P_i$  requests resource,  $P_i$  releases resources

$\mathbf{S}$  := set of possible states, e.g.  $\mathbf{S} =$  queue of waiting request commands, where request at head of queue is currently granted

$$\mathbf{e} : \mathbf{C} \times \mathbf{S} \rightarrow \mathbf{S}$$

$$\mathbf{e}(C, S) = S'$$

Synchronization achieved because all processes order commands according to their timestamps (using  $\implies$  relation), so each  $P_i$  uses same sequence of commands.

$P_i$  can execute command timestamped  $T$  where it has learned of all commands issued by all other processes with timestamps less than or equal to  $T$ .

## 52. VECTOR CLOCKS

A vector clock in a system of  $N$  processes (e.g.  $N = 2, N = 3$ ), is a vector of  $N$  integers.

Each process  $P_i$  maintains its own vector clock ( $V_i$  for a process  $P_i$ ) to timestamp local events.

Vector timestamps (vector of  $N$  integers) are sent with each message.

Rules for vector clocks:

- (1) vector is initialized to 0 at all processes:

$$(289) \quad V_i[j] = 0 \forall i, j = 1, \dots, N$$

- (2) Before a process  $P_i$  timestamps an event, it increments its element of the vector (index  $i$ ) in its local vector (vector of  $N$  integers):

$$(290) \quad V_i[i] = V_i[i] + 1$$

- (3) A message is sent from process  $P_i$  with  $V_i$  attached to the message.
- (4) When a process  $P_i$  receives a vector timestamp  $t$ , it compares the 2 vectors element by element, setting its local vector clock to the higher of the 2 values:

$$(291) \quad V_j[i] = \max(V_j[i], t[i]) \forall i = 1, \dots, N$$

Remember also to, upon the receiving event, increment the local vector of the receiving process, before the comparison.

We compare 2 vector timestamps by defining:

$$(292) \quad \begin{aligned} V = V' &\iff V[j] = V'[j] \forall i = 1, \dots, N \\ V \leq V' &\iff V[j] = V'[j] \forall i = 1, \dots, N \end{aligned}$$

For any 2 events  $e, e'$ , if  $e \rightarrow e'$ , then  $V(e) < V(e')$  which is the same as we get from Lamport's algorithm, Lmaport's timestamps.

But with vector clocks, we now have additional knowledge that if  $V(e) < V(e')$  then  $e \rightarrow e'$ . 2 events  $e, e'$  are concurrent if *neither*  $V(e) \leq V(e')$  nor  $V(e') \leq V(e)$ .

## 53. CLOCK SYNCHRONIZATION, LAMPORT TIMESTAMPS, VECTOR CLOCKS, REFERENCES

<https://www.cs.rutgers.edu/~pxk/rutgers/notes/content/08-logical-clocks-slides.pdf>

<https://towardsdatascience.com/understanding-lamport-timestamps-with-pythons-multiprocessing-library-12a642788>

Rutgers University – CS 417: Distributed Systems 2009 Paul Krzyzanowski, Lectures on distributed systems Clock Synchronization Paul Krzyzanowsk <https://www.cs.rutgers.edu/~pxk/rutgers/notes/content/08-clocks.pdf>

## Part 18. Flight Software

REFERENCES

[1] Gordon C. Oates. **Aerothermodynamics of Gas Turbine and Rocket Propulsion** 3rd Edition. AIAA; 3rd edition (January 1, 1997). ISBN-13: 978-1563472411

[2] Philip Hill and Carl Peterson. **Mechanics and Thermodynamics of Propulsion**. 2nd Ed. Pearson.

[3] George P. Sutton, Oscar Biblarz. **Rocket Propulsion Elements**, 7th Edition. Wiley, 2001.

[4] Ronald W. Humble, Gary N. Henry, Wiley J. Larson. **Space Propulsion Analysis and Design**. First Edition-Revised. The McGraw-Hill Companies, Inc. Primis Custom Publishing. 1995. ISBN 0-07-031320-2

[5] L D Landau, E.M. Lifshitz. **Mechanics: Volume 1** (Course of Theoretical Physics S) 3rd Edition. Butterworth-Heinemann; 3 edition (January 15, 1976). ISBN-10: 0750628960

[6] T. Frankel, **The Geometry of Physics**. Cambridge University Press, Third Edition, 2011.

[7] Ralph Baierlein. **Thermal Physics** Cambridge University Press (July 28, 1999), ISBN-13: 978-0521658386

[8] Charles Kittel, Herbert Kroemer, **Thermal Physics**, W. H. Freeman; Second Edition edition, 1980. ISBN-13: 978-0716710882

[9] Bernard F. Schutz, **Geometrical Methods of Mathematical Physics**, Cambridge University Press, 1980. ISBN-13: 978-0521298872

[10] Claus Borgnakke, Richard E. Sonntag. **Fundamentals of Thermodynamics**, 8th Edition, Wiley, (December 26, 2012). ISBN-13: 978-1118131992

[11] Herbert Goldstein, Charles P. Poole Jr., John L. Safko. **Classical Mechanics** (3rd Edition). Addison-Wesley; 3 edition (June 25, 2001). ISBN-13: 978-0201657029 <https://homepages.dias.ie/ydri/Goldstein.pdf>

[12] Joseph M. Powers. “Lecture Notes on Fundamentals of Combustion.” updated 30 March 2014, 2:12pm <http://www3.nd.edu/~powers/ame.60636/notes.pdf>

[13] Stephen Turns. **An Introduction to Combustion: Concepts and Applications** 3rd Edition. McGraw-Hill Education; 3 edition (January 24, 2011). The second edition, 2000, was used in these notes.

[14] L.D. Landau, E.M. Lifshitz. **Statistical Physics**, Third Edition, Part 1: Volume 5 (Course of Theoretical Physics, Volume 5). Butterworth-Heinemann; 3rd edition (January 15, 1980). ISBN-13: 978-0750633727

[15] Michel Le Bellac, Fabrice Mortessagne, G. George Batrouni. **Equilibrium and Non-Equilibrium Statistical Thermodynamics**. Cambridge University Press (May 3, 2004). ISBN-13: 978-0521821438

[16] John Lee, **Introduction to Smooth Manifolds** (Graduate Texts in Mathematics, Vol. 218), 2nd edition, Springer, 2012, ISBN-13: 978-1441999818

[17] Ovidiu Calin, Der-Chen Chang. **Geometric Mechanics on Riemannian Manifolds: Applications to Partial Differential Equations** (Applied and Numerical Harmonic Analysis). Birkhäuser. 2005. ISBN-13: 978-0817643546

[18] D. G. Goodwin, CANTERA, Division of Engineering and Applied Science, California Institute of Technology, [www.cantera.org](http://www.cantera.org).

[19] Philip A. Thompson. **Compressible-Fluid Dynamics**. 1988.

[20] Manuel Martinez-Sanchez. 16.512 Rocket Propulsion, Fall 2005. (Massachusetts Institute of Technology: MIT OpenCourseWare), <http://ocw.mit.edu> (Accessed 10 Feb, 2016). License: [Creative Commons BY-NC-SA](#)

[21] Joel H. Ferziger, Milovan Peric. **Computational Methods for Fluid Dynamics**, 3rd Edition. Springer. 2013.

[22] Tuncer Cebeci, Jian P. Shao, Fassi Kafyeke, Eric Laurendeau. **Computational Fluid Dynamics for Engineers**. Horizons Publishing Inc., 2005. ISBN 0-9766545-0-4 Horizons Pbulishing Inc., Long Beach, ISBN 3-540-24451-4 Springer Berlin Heidelberg, New York.

[23] Edited by Vigor Yang, Mohammed Habiballah, James Hulka, Michael Popp. **Liquid Rocket Thrust Chambers: Aspects of Modeling, Analysis, and Design**. Volume 200, Progress in Astronautics and Aeronautics. Paul Zarchan, Editor-in-Chief. 2004. ISBN 1-56347-223-6

[24] S.P. Novikov. “The Hamiltonian formalism and a many-valued analogue of Morse theory.” <http://www.mi.ras.ru/~snovikov/74.pdf>

[25] Yuanxun Bill Bao and Justin Meskas. “Lattice Boltzmann Method for Fluid Simulations.” <http://www.cims.nyu.edu/~billbao/report930.pdf>

[26] Kyle E. Niemeyer, Chih-Jen Sung. *Accelerating reactive-flow simulations using graphics processing units*. **American Institute of Aeronautics and Astronautics, Inc. (AIAA)**. *51st AIAA Aerospace Sciences Meeting including the New Horizons Forum and Aerospace Exposition; 07 - 10 January 2013, Grapevine (Dallas/Ft. Worth Region), Texas*. **Niemeyer-Sung-ASM\_2013.pdf**

[27] Robert G. Jahn. **Physics of Electric Propulsion** (Dover Books on Physics). Dover Publications (May 26, 2006).

[28] Edward M. Purcell. **Electricity and Magnetism** (Berkeley Physics Course, Vol. 2) 2nd. Edition. McGraw-Hill Science/Engineering/Math; 2 edition (August 1, 1984) ISBN-13: 978-0070049086

There’s a 3rd. edition, 2013, and another co-author, David Morin, but I don’t have access to a copy.

[29] Richard Fitzpatrick. **Plasma Physics** 2014 [https://www.cfa.harvard.edu/~scanmer/Ay253/LecNotes/fitzpatrick\\_plasma\\_physics.pdf](https://www.cfa.harvard.edu/~scanmer/Ay253/LecNotes/fitzpatrick_plasma_physics.pdf)

[30] L. P. Pitaevskii, E.M. Lifshitz. **Physical Kinetics: Volume 10** (Course of Theoretical Physics S) 1st Edition. Butterworth-Heinemann; 1 edition (January 15, 1981). ISBN-13: 978-0750626354

[31] Jonas Tölke. *Implementation of a Lattice Boltzmann kernel using the Compute Unified Device Architecture developed by nVIDIA*. Comput Visual Sci. DOI 10.1007/s00791-008-0120-2 [http://moodle.epfl.ch/pluginfile.php/952831/mod\\_resource/content/0/toelked2q9.pdf](http://moodle.epfl.ch/pluginfile.php/952831/mod_resource/content/0/toelked2q9.pdf)

[32] Y.H. Qian, D. D’Humières and P. Lallemand. Lattice BGK Models for Navier-Stokes Equation. **Europhysics Letters**. *Europhys. Lett.*, **17** (6), pp. 479-484 (1992). 1 February bis 1992

[33] Eric Sonnendrücker. *Advanced Finite Element Methods*. Lecture notes. Wintersemester 2016/2017. January 30, 2017. Eric Sonnendrücker, *Max-Planck-Institut für Plasmaphysik und Zentrum Mathematik, TU München*

[34] David Darmofal. \*16.901 Computational Methods in Aerospace Engineering, Spring 2005.\* (Massachusetts Institute of Technology: MIT OpenCourseWare), <http://ocw.mit.edu> (Accessed 12 Jun, 2016). License: [Creative Commons BY-NC-SA](#)

[35] Roger R. Bate, Donald D. Mueller, Jerry E. White. **Fundamentals of Astrodynamics** (Dover Books on Aeronautical Engineering) Revised ed. Edition. Dover Publications; Revised ed. edition (June 1, 1971). ISBN-10: 0486600610

[36] William H. Press, Saul A. Teukolsky, William T. Vetterling, Brian P. Flannery. **Numerical Recipes: The Art of Scientific Computing**. 3rd Edition. Cambridge University Press. ISBN-10: 0521880688

[37] E. Hairer, S.P. Nørsett, G. Wanner. **Solving Ordinary Differential Equations I: Nonstiff Problems**. Second Revised Edition. Springer. 1993.

[38] E. Hairer, G. Wanner. **Solving Ordinary Differential Equations II: Stiff and Differential-Algebraic Problems**. Second Revised Edition. Springer. 2010.

[39] Lawrence F. Shampine. ”Some Practical Runge-Kutta Formulas.” **Mathematics of Computuaton**. Volume 46. Number 173. Jan. 1986, Pages 135-150.

[40] Daniel Gaspar, Jack Stouffer. **Mastering Flask Web Development: Build enterprise-grade, scalable Python web applications**, 2nd Edition Kindle Edition. 2018.

[41] Howard D. Curtis. **Orbital Mechanics for Engineering Students**. Butterworth-Heinemann. Fourth Edition. 2020.

[42] Bong Wie. **Space Vehicle Dynamics and Control** (AIAA Education) 2nd Edition. (August 28, 2008)

[43] Dr. Murat Arcak. *Signals and Systems*. Fall 2019 (Fa19). UC Berkeley. <https://inst.eecs.berkeley.edu/~ee120/fa19/>

[44] Alan V. Oppenheim and Alan S. Willsky. **Signals and Systems**. Prentice-Hall, 2nd ed., 1997. ISBN-13: 978-0138147570

[45] O. L. R. Jacobs. **Introduction to Control Theory**. 2nd Edition. Oxford Universty Press. 1994.

[46] [Brian Douglas](#). **Discrete Control**. <https://youtube.com/playlist?list=PLUMWjy5jgHKOMLv6Ksf-NHi7Ur8NRNU4Z>

[47] Katsuhiko Ogata. **Modern Control Engineering**. Pearson; 5th edition (August 25, 2009).

[48] Daniel J. Block, Karl J. Åström, Mark W. Spong. **The Reaction Wheel Pendulum**. Springer Cham. 2007.

[49] Leslie Lamport. *Time, Clocks and the Ordering of Events in a Distributed System*. **Communications of the ACM** **21**, 7 (July 1978), 558-565. <https://lamport.azurewebsites.net/pubs/time-clocks.pdf>

There’s an 8th edition of Biblarz and Sutton [3] for 2010 that I would like to have. If you find any of this material useful or if you’d like to help, email me or visit my Open/Tilt page [ernestyalumni.tilt.com](http://ernestyalumni.tilt.com) and donate to the crowdfunding campaign or click on the PayPal donate button.